

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации

# Лекция 2. Вероятностная постановка задачи распознавания образов. Обобщенные линейные модели

Д. П. Ветров<sup>1</sup>    Д. А. Кропотов<sup>2</sup>

<sup>1</sup>МГУ, ВМиК, каф. ММП

<sup>2</sup>ВЦ РАН

Спецкурс «Байесовские методы машинного обучения»

# План лекции

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации

## Ликбез

Нормальное распределение  
Решение нерешаемых СЛАУ

## Статистическая постановка задачи машинного обучения

Вероятностное описание  
Байесовские решающие правила

## Линейная регрессия

Классическая линейная регрессия  
Метод наименьших квадратов  
Вероятностная постановка задачи

## Задача классификации

Логистическая регрессия  
Метод IRLS

# План лекции

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез  
Нормальное  
распределение  
Решение  
нерешаемых  
СЛАУ

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации

## Ликбез

Нормальное распределение

Решение нерешаемых СЛАУ

Статистическая постановка задачи машинного обучения

Вероятностное описание

Байесовские решающие правила

Линейная регрессия

Классическая линейная регрессия

Метод наименьших квадратов

Вероятностная постановка задачи

Задача классификации

Логистическая регрессия

Метод IRLS

# Нормальное распределение

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Нормальное  
распределение  
Решение  
нерешаемых  
СЛАУ

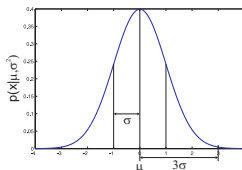
Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации

- Нормальное распределение играет важнейшую роль в математической статистике

$$X \sim \mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$
$$\mu = \mathbb{E}X, \quad \sigma^2 = \mathbb{D}X \triangleq \mathbb{E}(X - \mathbb{E}X)^2$$



- Из центральной предельной теоремы следует, что сумма независимых случайных величин с ограниченной дисперсией стремится к нормальному распределению
- На практике, многие случайные величины можно считать приближенно нормальными

# Многомерное нормальное распределение

Лекция 2.  
Вероятностная постановка задачи распознавания образов.  
Обобщенные линейные модели

Ветров,  
Кропотов

Ликбез

Нормальное распределение  
Решение нерешаемых СЛАУ

Статистическая постановка задачи машинного обучения

Линейная регрессия

Задача классификации

- Многомерное нормальное распределение имеет вид

$$X \sim \mathcal{N}(x|\mu, \Sigma) = \frac{1}{\sqrt{2\pi}^n \sqrt{\det \Sigma}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right),$$

где  $\mu = \mathbb{E}X$ ,  $\Sigma = \mathbb{E}(X - \mu)(X - \mu)^T$  — вектор математических ожиданий каждой из  $n$  компонент и матрица ковариаций соответственно

- Матрица ковариаций показывает, насколько сильно связаны (коррелируют) компоненты многомерного нормального распределения

$$\Sigma_{ij} = \mathbb{E}(X_i - \mu_i)(X_j - \mu_j) = \text{Cov}(X_i, X_j)$$

- Если мы поделим ковариацию на корень из произведений дисперсий, то получим коэффициент корреляции

$$\rho(X_i, X_j) \triangleq \frac{\text{Cov}(X_i, X_j)}{\sqrt{\mathbb{D}X_i \mathbb{D}X_j}} \in [-1, 1]$$

# Особенности нормального распределения

Лекция 2.  
Вероятностная постановка задачи распознавания образов.  
Обобщенные линейные модели

Ветров,  
Кропотов

Ликбез

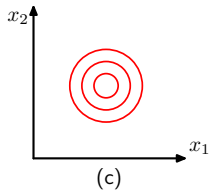
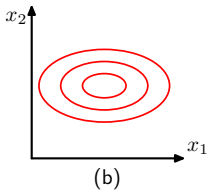
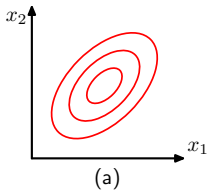
Нормальное распределение  
Решение нерешаемых СЛАУ

Статистическая постановка задачи машинного обучения

Линейная регрессия

Задача классификации

- Нормальное распределение **полностью задается** первыми двумя моментами (мат. ожидание и матрица ковариаций/дисперсия)
- Матрица ковариаций неотрицательно определена, причем на диагоналях стоят дисперсии соответствующих компонент
- Нормальное распределение имеет очень легкие хвосты: большие отклонения от мат. ожидания практически невозможны. Это обстоятельство нужно учитывать при приближении произвольных случайных величин нормальными



# План лекции

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез  
Нормальное  
распределение  
Решение  
нерешаемых  
СЛАУ

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации

## Ликбез

Нормальное распределение  
Решение нерешаемых СЛАУ

Статистическая постановка задачи машинного обучения  
Вероятностное описание  
Байесовские решающие правила

## Линейная регрессия

Классическая линейная регрессия  
Метод наименьших квадратов  
Вероятностная постановка задачи

## Задача классификации

Логистическая регрессия  
Метод IRLS

# Псевдообращение матриц

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез  
Нормальное  
распределение  
Решение  
нерешаемых  
СЛАУ

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации

- Предположим, нам необходимо решить СЛАУ вида  $Ax = b$
- Если бы матрица  $A$  была квадратной и невырожденной (число уравнений равно числу неизвестных и все уравнения линейно независимы), то решение задавалось бы формулой  $x = A^{-1}b$
- Предположим, что число уравнений больше числа неизвестных, т.е. матрица  $A$  прямоугольная. Домножим обе части уравнения на  $A^T$  слева

$$A^T Ax = A^T b$$

- В левой части теперь квадратная матрица и ее можно перенести в правую часть

$$x = (A^T A)^{-1} A^T b$$

- Операция  $(A^T A)^{-1} A^T$  называется псевдообращением матрицы  $A$ , а  $x$  – псевдорешением



# Нормальное псевдорешение

Лекция 2.

Вероятностная постановка задачи распознавания образов. Обобщенные линейные модели

Ветров, Кропотов

Ликбез  
Нормальное распределение  
Решение нерешаемых СЛАУ

Статистическая постановка задачи машинного обучения

Линейная регрессия

Задача классификации

- Если матрица  $A^T A$  вырождена, псевдорешений бесконечно много, причем найти их на компьютере нетривиально
- Для решения этой проблемы используется ридж-регуляризация матрицы  $A^T A$

$$A^T A + \lambda I,$$

где  $I$  – единичная матрица, а  $\lambda$  – коэффициент регуляризации. Такая матрица невырождена для любых  $\lambda > 0$

- Величина

$$\mathbf{x} = (A^T A + \lambda I)^{-1} A^T \mathbf{b}$$

называется нормальным псевдорешением. Оно всегда единственно и при небольших положительных  $\lambda$  определяет псевдорешение с наименьшей нормой

# Графическая иллюстрация

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

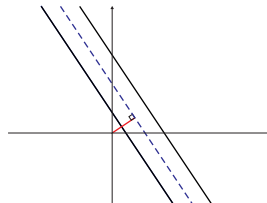
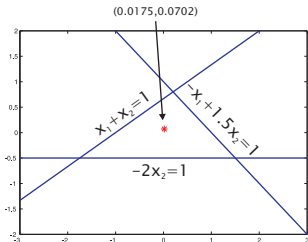
Ликбез  
Нормальное  
распределение  
Решение  
нерешаемых  
СЛАУ

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации

- Псевдорешение соответствует точке, минимизирующей невязку, а нормальное псевдорешение отвечает псевдорешению с наименьшей нормой



- Заметим, что псевдообратная матрица  $(A^T A)^{-1} A^T$  совпадает с обратной матрицей  $A^{-1}$  в случае невырожденных квадратных матриц

# План лекции

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Вероятностное  
описание

Байесовские  
решающие  
правила

Линейная  
регрессия

Задача  
классификации

## Ликбез

Нормальное распределение  
Решение нерешаемых СЛАУ

## Статистическая постановка задачи машинного обучения

Вероятностное описание  
Байесовские решающие правила

## Линейная регрессия

Классическая линейная регрессия  
Метод наименьших квадратов  
Вероятностная постановка задачи

## Задача классификации

Логистическая регрессия  
Метод IRLS

# Основные обозначения

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Вероятностное  
описание

Байесовские  
решающие  
правила

Линейная  
регрессия

Задача  
классификации

- В дальнейшем будут рассматриваться преимущественно задачи классификации и восстановления регрессии
- В этих задачах обучающая выборка представляет собой набор отдельных объектов  $X = \{\mathbf{x}_i\}_{i=1}^n$ , характеризующихся вектором вещественнозначных признаков  $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d})$
- Каждый объект также обладает скрытой переменной  $t \in \mathcal{T}$
- Предполагается, что существует зависимость между признаками объекта и значением скрытой переменной
- Для объектов обучающей выборки значение скрытой переменной известно  $\mathbf{t} = \{t_i\}_{i=1}^n$

# Статистическая постановка задачи

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Вероятностное  
описание

Байесовские  
решающие  
правила

Линейная  
регрессия

Задача  
классификации

- Каждый объект описывается парой  $(\mathbf{x}, t)$
- При статистической (вероятностной) постановке задачи машинного обучения предполагается, что обучающая выборка является набором независимых, одинаково распределенных случайных величин, взятых из некоторой генеральной совокупности
- В этом случае уместно говорить о плотности распределения объектов  $p(\mathbf{x}, t)$  и использовать вероятностные термины (математическое ожидание, дисперсия, правдоподобие) для описания и решения задачи
- Заметим, что это не единственная возможная постановка задачи машинного обучения

# Качество обучения

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Вероятностное  
описание

Байесовские  
решающие  
правила

Линейная  
регрессия

Задача  
классификации

- Качество обучения определяется точностью прогноза на генеральной совокупности
- Пусть  $S(t, \hat{t})$  – функция потерь, определяющая штраф за прогноз  $\hat{t}$  при истинном значении скрытой переменной  $t$
- Разумно ожидать, что минимум этой функции достигается при  $\hat{t} = t$
- Примерами могут служить  $S_r(t, \hat{t}) = (t - \hat{t})^2$  для задачи восстановления регрессии и  $S_c(t, \hat{t}) = I\{\hat{t} \neq t\}$  для задачи классификации

# Абсолютный критерий качества

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Вероятностное  
описание

Байесовские  
решающие  
правила

Линейная  
регрессия

Задача  
классификации

- Если бы функция  $p(\mathbf{x}, t)$  была известна, задачи машинного обучения не существовало
- В самом деле абсолютным критерием качества обучения является мат. ожидание функции потерь, взятое по генеральной совокупности

$$\mathbb{E}S(t, \hat{t}) = \int S(t, \hat{t}(\mathbf{x}))p(\mathbf{x}, t)d\mathbf{x}dt \rightarrow \min,$$

где  $\hat{t}(\mathbf{x})$  – решающее правило, возвращающее величину прогноза для вектора признаков  $\mathbf{x}$

- Вместо методов машинного обучения сейчас бы активно развивались методы оптимизации и взятия интегралов от функции потерь :)
- К сожалению (а может, к счастью), распределение объектов генеральной совокупности неизвестно, поэтому абсолютный критерий качества обучения не может быть подсчитан

# План лекции

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Вероятностное  
описание

Байесовские  
решающие  
правила

Линейная  
регрессия

Задача  
классификации

## Ликбез

Нормальное распределение  
Решение нерешаемых СЛАУ

## Статистическая постановка задачи машинного обучения

Вероятностное описание  
Байесовские решающие правила

## Линейная регрессия

Классическая линейная регрессия  
Метод наименьших квадратов  
Вероятностная постановка задачи

## Задача классификации

Логистическая регрессия  
Метод IRLS



# Идеальный классификатор

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Вероятностное  
описание

Байесовские  
решающие  
правила

Линейная  
регрессия

Задача  
классификации

- Итак, одна из основных задач теории машинного обучения — это разработка способов косвенного оценивания качества решающего правила и выработка новых критериев для оптимизации в ходе обучения
- Рассмотрим задачу классификации с функцией потерь вида  $S_c(t, \hat{t}) = I\{\hat{t} \neq t\}$  и гипотетический классификатор  $t_B(\mathbf{x}) = \arg \max_{t \in \mathcal{T}} p(\mathbf{x}, t) = \arg \max_{t \in \mathcal{T}} p(t|\mathbf{x})$
- Справедлива следующая цепочка неравенств

$$\begin{aligned}\mathbb{E}S(t, \hat{t}) &= \int \int S(t, \hat{t}(\mathbf{x}))p(\mathbf{x}, t)d\mathbf{x}dt = \\ &= \sum_{s=1}^l \int S(s, \hat{t}(\mathbf{x}))p(\mathbf{x}, s)d\mathbf{x} = 1 - \int p(\mathbf{x}, \hat{t}(\mathbf{x}))d\mathbf{x} \geq \\ &\geq 1 - \int \max_t p(\mathbf{x}, t)d\mathbf{x} = 1 - \int p(\mathbf{x}, t_B(\mathbf{x}))d\mathbf{x} = \mathbb{E}S(t, t_B)\end{aligned}$$

# Идеальная регрессия

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения  
Вероятностное  
описание

Байесовские  
решающие  
правила

Линейная  
регрессия

Задача  
классификации

- Рассмотрим задачу восстановления регрессии с квадратичной функцией потерь вида  $S_r(t, \hat{t}) = (t - \hat{t})^2$  и гипотетическое решающее правило

$$t_B(\mathbf{x}) = \mathbb{E}_{t|\mathbf{x}} t = \int t p(t|\mathbf{x}) dt$$

- Справедлива следующая цепочка неравенств

$$\begin{aligned} \mathbb{E} S(t, \hat{t}) &= \int \int S(t, \hat{t}(\mathbf{x})) p(\mathbf{x}, t) d\mathbf{x} dt = \\ &= \int \int (t - \hat{t}(\mathbf{x}))^2 p(\mathbf{x}, t) d\mathbf{x} dt = \int \int ((t - \mathbb{E}t) + (\mathbb{E}t - \hat{t}(\mathbf{x})))^2 p(\mathbf{x}, t) d\mathbf{x} dt = \\ &= \int \int (t - \mathbb{E}t)^2 p(\mathbf{x}, t) d\mathbf{x} dt + 2 \int \int (t - \mathbb{E}t)(\mathbb{E}t - \hat{t}(\mathbf{x})) p(\mathbf{x}, t) d\mathbf{x} dt + \\ &\quad + \int \int (\mathbb{E}t - \hat{t}(\mathbf{x}))^2 p(\mathbf{x}, t) d\mathbf{x} dt \geq \\ &\geq \int \int (t - \mathbb{E}t)^2 p(t|\mathbf{x}) p(\mathbf{x}) dt d\mathbf{x} = \mathbb{E} S(t, t_B(\mathbf{x})) \end{aligned}$$

# Особенности байесовских решающих правил

Лекция 2.

Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения  
Вероятностное  
описание

Байесовские  
решающие  
правила

Линейная  
регрессия

Задача  
классификации

- Таким образом, знание распределения объектов генеральной совокупности приводит к получению оптимальных решающих правил **в явной форме**
- Такой оптимальные решающие правила называются байесовскими
- Если бы удалось с высокой точностью оценить значение условной плотности  $p(t|\mathbf{x})$  для всех  $\mathbf{x}$  и  $t$ , обе основные задачи машинного обучения можно было считать решенными
- На этом основан один из существующих подходов к машинному обучению

# План лекции

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Классическая  
линейная  
регрессия

Метод  
наименьших  
квадратов  
Вероятностная  
постановка  
задачи

Задача  
классификации

## Ликбез

Нормальное распределение  
Решение нерешаемых СЛАУ

## Статистическая постановка задачи машинного обучения

Вероятностное описание  
Байесовские решающие правила

## Линейная регрессия

Классическая линейная регрессия  
Метод наименьших квадратов  
Вероятностная постановка задачи

## Задача классификации

Логистическая регрессия  
Метод IRLS

# Задача восстановления регрессии

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Классическая  
линейная  
регрессия

Метод  
наименьших  
квадратов

Вероятностная  
постановка  
задачи

Задача  
классификации

- Задача восстановления регрессии предполагает наличие связи между наблюдаемыми признаками  $\mathbf{x}$  и непрерывной переменной  $t$
- В отличие от задачи интерполяции допускаются отклонения решающего правила от правильных ответов на объектах обучающей выборки
- Уравнение регрессии  $y(\mathbf{x}, \mathbf{w})$  ищется в некотором параметрическом виде путем нахождения наилучшего значения вектора весов

$$\mathbf{w}_* = \arg \max_{\mathbf{w}} F(X, t, \mathbf{w})$$

# Линейная регрессия

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Классическая  
линейная  
регрессия

Метод  
наименьших  
квадратов

Вероятностная  
постановка  
задачи

Задача  
классификации

- Наиболее простой и изученной является линейная регрессия
- Главная особенность: настраиваемые параметры входят в решающее правило **линейно**
- Заметим, что линейная регрессия не обязана быть линейной по признакам
- Общее уравнение регрессии имеет вид

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=1}^m w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

# Особенность выбора базисных функций

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Классическая  
линейная  
регрессия

Метод  
наименьших  
квадратов  
Вероятностная  
постановка  
задачи

Задача  
классификации

- Общего метода выбора базисных функций  $\phi_j(\mathbf{x})$  — не существует
- Обычно они подбираются из априорных соображений (например, если мы пытаемся восстановить какой-то периодический сигнал, разумно взять функции тригонометрического ряда) или путем использования некоторых «универсальных» базисных функций
- Наиболее распространенными базисными функциями являются
  - $\phi(\mathbf{x}) = x_k$
  - $\phi(\mathbf{x}) = x_{k_1} x_{k_2} \dots x_{k_l}$
  - $\phi(\mathbf{x}) = \exp(-\gamma \|\mathbf{x} - \mathbf{x}_0\|^p)$ ,  $\gamma, p > 0$ .
- Метод построения линейной регрессии (настройки весов  $\mathbf{w}$ ) **не зависит** от выбора базисных функций

# Формализация задачи

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Классическая  
линейная  
регрессия

Метод  
наименьших  
квадратов

Вероятностная  
постановка  
задачи

Задача  
классификации

- Пусть  $S(t, \hat{t})$  — функция потерь от ошибки в определении регрессионной переменной  $t$
- Необходимо минимизировать потери от ошибок на генеральной совокупности

$$\mathbb{E}S(t, y(\mathbf{x}, \mathbf{w})) = \int \int S(t, y(\mathbf{x}, \mathbf{w}))p(\mathbf{x}, t)d\mathbf{x}dt \rightarrow \min_{\mathbf{w}}$$

- Дальнейшие рассуждения зависят от вида функции потерь
- Во многих случаях даже не нужно восстанавливать полностью условное распределение  $p(t|\mathbf{x})$



# Важная теорема

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Классическая  
линейная  
регрессия

Метод  
наименьших  
квадратов

Вероятностная  
постановка  
задачи

Задача  
классификации

- Теорема. Пусть функция потерь имеет вид
  - $S(t, \hat{t}) = (t - \hat{t})^2$  — «Потери старушки»;
  - $S(t, \hat{t}) = |t - \hat{t}|$  — «Потери олигарха»;
  - $S(t, \hat{t}) = \delta^{-1}(t - \hat{t})$  — «Потери инвалида».

Тогда величиной, минимизирующей функцию  $\mathbb{E}S(t, y(\mathbf{x}, \mathbf{w}))$ , является следующая

- $y(\mathbf{x}) = \mathbb{E}p(t|\mathbf{x})$ ;
  - $y(\mathbf{x}) = \text{med } p(t|\mathbf{x})$ ;
  - $y(\mathbf{x}) = \text{mod } p(t|\mathbf{x}) = \arg \max_t p(t|\mathbf{x})$ .
- В зависимости от выбранной системы предпочтений, мы будем пытаться оценивать тот или иной функционал от апостериорного распределения **вместо того, чтобы оценивать его самого**

# План лекции

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Классическая  
линейная  
регрессия

Метод  
наименьших  
квадратов

Вероятностная  
постановка  
задачи

Задача  
классификации

## Ликбез

Нормальное распределение  
Решение нерешаемых СЛАУ

## Статистическая постановка задачи машинного обучения

Вероятностное описание  
Байесовские решающие правила

## Линейная регрессия

Классическая линейная регрессия  
Метод наименьших квадратов  
Вероятностная постановка задачи

## Задача классификации

Логистическая регрессия  
Метод IRLS

# Минимизация невязки

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Классическая  
линейная  
регрессия

Метод  
наименьших  
квадратов

Вероятностная  
постановка  
задачи

Задача  
классификации

- Наиболее часто используемой функцией потерь является квадратичная  $S(t, \hat{t}) = (t - \hat{t})^2$
- Значение регрессионной функции на обучающей выборке в матричном виде может быть записано как  $\mathbf{y} = \Phi \mathbf{w}$ , где  $\Phi = (\phi_{ij}) = (\phi_j(\mathbf{x}_i)) \in \mathbb{R}^{n \times m}$
- Таким образом, приходим к следующей задаче

$$\|\mathbf{y} - \mathbf{t}\|^2 = \|\Phi \mathbf{w} - \mathbf{t}\|^2 \rightarrow \min_{\mathbf{w}}$$

Взяв производную по  $\mathbf{w}$  и приравняв ее к нулю, получаем

$$\begin{aligned} \frac{\partial \|\Phi \mathbf{w} - \mathbf{t}\|^2}{\partial \mathbf{w}} &= \frac{\partial [\mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{w}^T \Phi^T \mathbf{t} + \mathbf{t}^T \mathbf{t}]}{\partial \mathbf{w}} = \\ &= 2\Phi^T \Phi \mathbf{w} - 2\Phi^T \mathbf{t} = 0 \\ \mathbf{w} &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \end{aligned}$$

# Регуляризация задачи

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия  
Классическая  
линейная  
регрессия

Метод  
наименьших  
квадратов

Вероятностная  
постановка  
задачи

Задача  
классификации

- Заметим, что формула для весов линейной регрессии представляет собой псевдорешение уравнения  $\Phi \mathbf{w} = \mathbf{t}$
- Матрица  $\Phi^T \Phi \in \mathbb{R}^{m \times m}$  вырождена (Упр.) при  $m > n$
- Регуляризуя вырожденную матрицу, получаем

$$\mathbf{w} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{t}$$

- Отсюда формула для прогноза объектов обучающей выборки по их правильным значениям

$$\hat{\mathbf{t}} = \mathbf{y} = \Phi (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{t} = H \mathbf{t}$$

С историческим обозначением прогноза — навешиванием шляпки связано неформальное название матрицы  $H$ , по-английски звучащее как hat-matrix

# Особенности квадратичной функции потерь

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия  
Классическая  
линейная  
регрессия

Метод  
наименьших  
квадратов

Вероятностная  
постановка  
задачи

Задача  
классификации

- Достоинства
  - Квадратичная функция потерь гладкая (непрерывная и дифференцируемая)
  - Решение может быть получено в явном виде
  - Существует простая вероятностная интерпретация прогноза и функции потерь
- Недостатки
  - Решение неустойчиво (не робастно) относительно даже малого количества выбросов. Это связано с быстрым возрастанием квадратичной функции потерь при больших отклонениях от нуля
  - Квадратичная функция неприменима к задачам классификации

# План лекции

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Классическая  
линейная  
регрессия  
Метод  
наименьших  
квадратов

Вероятностная  
постановка  
задачи

Задача  
классификации

## Ликбез

Нормальное распределение  
Решение нерешаемых СЛАУ

## Статистическая постановка задачи машинного обучения

Вероятностное описание  
Байесовские решающие правила

## Линейная регрессия

Классическая линейная регрессия  
Метод наименьших квадратов  
Вероятностная постановка задачи

## Задача классификации

Логистическая регрессия  
Метод IRLS

# Нормальное распределение ошибок

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия  
Классическая  
линейная  
регрессия  
Метод  
наименьших  
квадратов

Вероятностная  
постановка  
задачи

Задача  
классификации

- Рассмотрим вероятностную постановку задачи восстановления регрессии. Регрессионная переменная  $t$  — случайная величина с плотностью распределения  $p(t|\mathbf{x})$
- В большинстве случаев предполагается, что  $t$  распределена нормально относительно некоторого мат. ожидания  $y(\mathbf{x})$ , определяемого точкой  $\mathbf{x}$

$$t = y(\mathbf{x}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(\varepsilon|0, \sigma^2)$$

- Необходимо найти функцию  $y(\mathbf{x})$ , которую мы можем отождествить с уравнением регрессии
- Предположение о нормальном распределении отклонений можно обосновать ссылкой на центральную предельную теорему

# Метод максимального правдоподобия для регрессии

Лекция 2.  
Вероятностная постановка задачи распознавания образов.  
Обобщенные линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая постановка задачи машинного обучения

Линейная регрессия  
Классическая линейная регрессия  
Метод наименьших квадратов  
Вероятностная постановка задачи

Задача классификации

- Используем ММП (не путать с одноименной кафедрой) для поиска  $y(\mathbf{x})$
- Правдоподобие задается следующей формулой

$$p(\mathbf{t}|\mathbf{y}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t_i - y_i)^2}{2\sigma^2}\right) \rightarrow \max$$

- Взяв логарифм и отбросив члены, не влияющие на положение максимума, получим

$$\sum_{i=1}^n (t_i - y_i)^2 = \sum_{i=1}^n (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 \rightarrow \min_{\mathbf{w}}$$

- Таким образом, применение метода максимального правдоподобия **в предположении о нормальности отклонений** эквивалентно методу наименьших квадратов



# Вероятностный смысл регуляризации

Лекция 2.  
Вероятностная постановка задачи распознавания образов.  
Обобщенные линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая постановка задачи машинного обучения

Линейная регрессия  
Классическая линейная регрессия  
Метод наименьших квадратов  
Вероятностная постановка задачи

Задача классификации

- Теперь будем максимизировать не правдоподобие, а апостериорную вероятность
- По формуле условной вероятности

$$p(\mathbf{w}|\mathbf{t}, X) = \frac{p(\mathbf{t}|X, \mathbf{w})p(\mathbf{w})}{p(\mathbf{t}, X)} \rightarrow \max_{\mathbf{w}}$$

знаменатель не зависит от  $\mathbf{w}$ , поэтому им можно пренебречь

- Пусть  $p(\mathbf{w}) \sim \mathcal{N}\left(\mathbf{w} \mid \mathbf{0}, \left(\frac{\sigma^2}{\lambda}\right) I\right)$ . Тогда

$$p(\mathbf{w}|\mathbf{t}, X) \propto \frac{\lambda^{m/2}}{(\sqrt{2\pi}\sigma)^{m+n}} \exp\left(-\frac{1}{2}\left(\sigma^{-2}\|\Phi\mathbf{w} - \mathbf{t}\|^2 + \frac{\lambda}{\sigma^2}\|\mathbf{w}\|^2\right)\right)$$

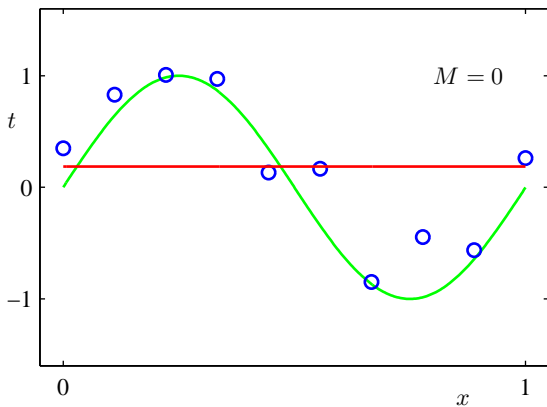
- Логарифмируя и приравнявая производную по  $\mathbf{w}$  к нулю, получаем

$$\mathbf{w} = (\Phi^T\Phi + \lambda I)^{-1}\Phi^T\mathbf{t}$$

- Регуляризация эквивалентна введению априорного распределения, поощряющего небольшие веса

# Зачем нужна регуляризация весов

Рассмотрим задачу восстановления регрессии с полиномиальными базисными функциями:  $x \in \mathbb{R}$ ,  $\phi_j(x) = x^j$ ,  $j = 0, \dots, M$



Лекция 2.

Вероятностная постановка задачи распознавания образов. Обобщенные линейные модели

Ветров, Кропотов

Ликбез

Статистическая постановка задачи машинного обучения

Линейная регрессия

Классическая линейная регрессия

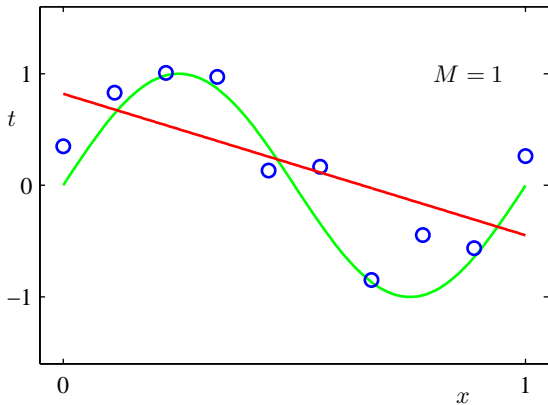
Метод наименьших квадратов

Вероятностная постановка задачи

Задача классификации

# Зачем нужна регуляризация весов

Рассмотрим задачу восстановления регрессии с полиномиальными базисными функциями:  $x \in \mathbb{R}$ ,  $\phi_j(x) = x^j$ ,  $j = 0, \dots, M$



Лекция 2.

Вероятностная постановка задачи распознавания образов. Обобщенные линейные модели

Ветров, Кропотов

Ликбез

Статистическая постановка задачи машинного обучения

Линейная регрессия

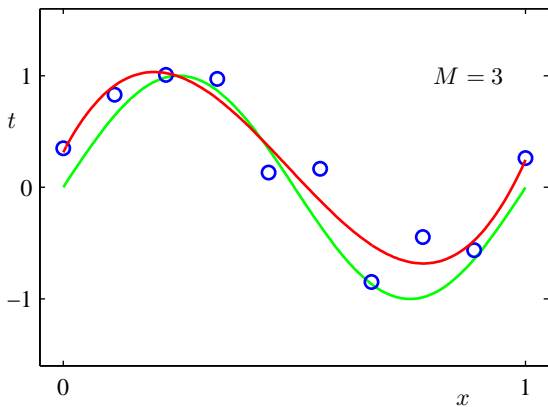
Классическая линейная регрессия  
Метод наименьших квадратов

Вероятностная постановка задачи

Задача классификации

# Зачем нужна регуляризация весов

Рассмотрим задачу восстановления регрессии с полиномиальными базисными функциями:  $x \in \mathbb{R}$ ,  $\phi_j(x) = x^j$ ,  $j = 0, \dots, M$



Лекция 2.

Вероятностная постановка задачи распознавания образов.

Обобщенные линейные модели

Ветров, Кропотов

Ликбез

Статистическая постановка задачи машинного обучения

Линейная регрессия

Классическая линейная регрессия

Метод наименьших квадратов

Вероятностная постановка задачи

Задача классификации

# Зачем нужна регуляризация весов

Лекция 2.  
Вероятностная постановка задачи распознавания образов.  
Обобщенные линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая постановка задачи машинного обучения

Линейная регрессия

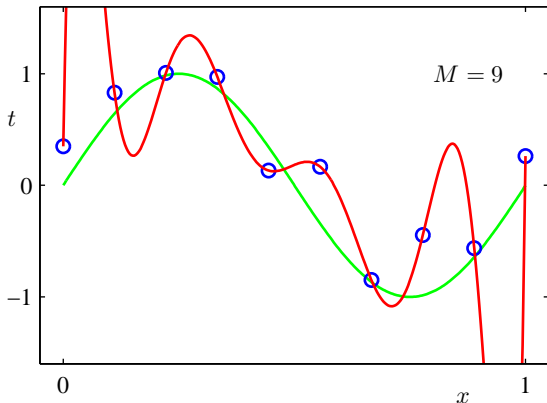
Классическая линейная регрессия

Метод наименьших квадратов

Вероятностная постановка задачи

Задача классификации

Рассмотрим задачу восстановления регрессии с полиномиальными базисными функциями:  $x \in \mathbb{R}$ ,  $\phi_j(x) = x^j$ ,  $j = 0, \dots, M$



# Значения наиболее правдоподобных весов

weight	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0$	0.19	0.82	0.31	0.35
$w_1$		-1.27	7.99	232.37
$w_2$			-25.43	-5321.83
$w_3$			17.37	48568.31
$w_4$				-231639.30
$w_5$				640042.26
$w_6$				-1061800.52
$w_7$				1042400.18
$w_8$				-557682.99
$w_9$				125201.43

Таблица: Значения наиболее правдоподобных весов в зависимости от степени полинома. С увеличением степени, абсолютные значения весов быстро растут

Лекция 2.  
Вероятностная постановка задачи распознавания образов.  
Обобщенные линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая постановка задачи машинного обучения

Линейная регрессия  
Классическая линейная регрессия  
Метод наименьших квадратов

Вероятностная постановка задачи

Задача классификации

# План лекции

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации

Логистическая  
регрессия  
Метод IRLS

## Ликбез

Нормальное распределение  
Решение нерешаемых СЛАУ

## Статистическая постановка задачи машинного обучения

Вероятностное описание  
Байесовские решающие правила

## Линейная регрессия

Классическая линейная регрессия  
Метод наименьших квадратов  
Вероятностная постановка задачи

## Задача классификации

Логистическая регрессия  
Метод IRLS

# Особенности задачи классификации

Лекция 2.

Вероятностная  
постановка  
задачи  
распознавания  
образов.

Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации  
Логистическая  
регрессия  
Метод IRLS

- Рассмотрим задачу классификации на два класса  $t \in \{-1, +1\}$
- Ее можно свести к задаче регрессии, например, следующим образом

$$\hat{t}(\mathbf{x}) = \text{sign}(y(\mathbf{x})) = \text{sign} \sum_{j=1}^m w_j \phi_j(\mathbf{x})$$

- Возникает вопрос: что использовать в качестве значений регрессионной переменной на этапе обучения?
- Наиболее распространенный подход заключается в использовании значения  $+\infty$  для  $t = +1$  и  $-\infty$  для  $t = -1$
- Геометрический смысл: чем дальше от нуля значение  $y(\mathbf{x})$ , тем увереннее мы в классификации объекта  $\mathbf{x}$



# Правдоподобие правильной классификации

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации  
Логистическая  
регрессия  
Метод IRLS

- Метод наименьших квадратов, очевидно, неприменим при таком подходе
- Воспользуемся вероятностной постановкой для выписывания функционала качества
- Определим правдоподобие классификации следующим образом

$$p(t|\mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-t\mathbf{y}(\mathbf{x}))}$$

- Это логистическая функция. Легко показать, что  $\sum_i p(t|\mathbf{x}, \mathbf{w}) = 1$  и  $p(t|\mathbf{x}, \mathbf{w}) > 0$ , а, значит, она является функцией правдоподобия

# Функционал качества в логистической регрессии

Лекция 2.  
Вероятностная постановка задачи распознавания образов.  
Обобщенные линейные модели

Ветров,  
Кропотов

Ликбез

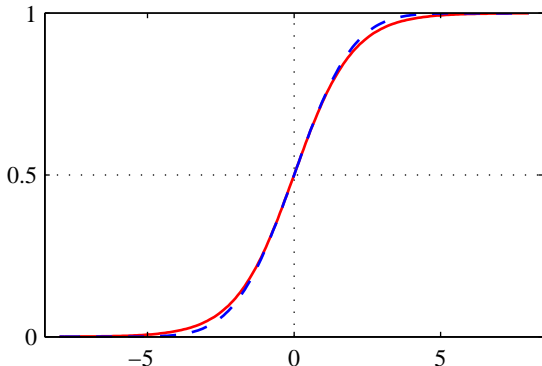
Статистическая постановка задачи машинного обучения

Линейная регрессия

Задача классификации

Логистическая регрессия

Метод IRLS



- Правдоподобие правильной классификации всей выборки имеет вид

$$p(t|X, \mathbf{w}) = \prod_{i=1}^n p(t_i|\mathbf{x}_i, \mathbf{w}) = \prod_{i=1}^n \frac{1}{1 + \exp\left(-t_i \sum_{j=1}^m w_j \phi_j(\mathbf{x}_i)\right)}$$

# План лекции

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации  
Логистическая  
регрессия  
Метод IRLS

## Ликбез

Нормальное распределение  
Решение нерешаемых СЛАУ

## Статистическая постановка задачи машинного обучения

Вероятностное описание  
Байесовские решающие правила

## Линейная регрессия

Классическая линейная регрессия  
Метод наименьших квадратов  
Вероятностная постановка задачи

## Задача классификации

Логистическая регрессия  
Метод IRLS

# Особенности функции правдоподобия классификации

Лекция 2.  
Вероятностная постановка задачи распознавания образов.  
Обобщенные линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая постановка задачи машинного обучения

Линейная регрессия

Задача классификации  
Логистическая регрессия  
Метод IRLS

- Приравнивание градиента логарифма правдоподобия к нулю приводит к трансцендентным уравнениям, которые неразрешимы аналитически
- Легко показать, что гессиан логарифма правдоподобия неположительно определен

$$\frac{\partial^2 \log p(\mathbf{t}|\mathbf{x}, \mathbf{w})}{\partial \mathbf{w}^2} \leq 0$$

- Это означает, что логарифм функции правдоподобия является вогнутым.
- Логарифм правдоподобия обучающей выборки  $L(\mathbf{w}) = \log p(\mathbf{t}|\mathbf{X}, \mathbf{w})$ , являющийся суммой вогнутых функций, также вогнут, а, значит, имеет **единственный максимум**

# Метод оптимизации Ньютона

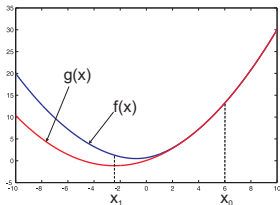
Основная идея метода Ньютона — это приближение в заданной точке оптимизируемой функции параболой и выбор минимума этой параболы в качестве следующей точки итерационного процесса:

$$f(\mathbf{x}) \rightarrow \min_{\mathbf{w}}$$

$$f(\mathbf{x}) \simeq g(\mathbf{x}) = f(\mathbf{x}_0) + (\nabla f(\mathbf{x}_0))^T(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T(\nabla \nabla f(\mathbf{x}_0))(\mathbf{x} - \mathbf{x}_0)$$

$$\nabla g(\mathbf{x}_*) = \nabla f(\mathbf{x}_0) + (\nabla \nabla f(\mathbf{x}_0))(\mathbf{x}_* - \mathbf{x}_0) = 0 \Rightarrow \mathbf{x}_* = \mathbf{x}_0 - (\nabla \nabla f(\mathbf{x}_0))^{-1}(\nabla f(\mathbf{x}_0))$$

Пример. Функция  $f(x) = \log(1 + \exp(x)) + \frac{x^2}{5}$ .  
 $x_0 = 6$ ,  $x_1 = -2.4418$ .



Лекция 2.

Вероятностная постановка задачи распознавания образов.

Обобщенные линейные модели

Ветров, Кропотов

Ликбез

Статистическая постановка задачи машинного обучения

Линейная регрессия

Задача классификации  
Логистическая регрессия  
Метод IRLS

# Итеративная минимизация логарифма правдоподобия

Лекция 2.

Вероятностная постановка задачи распознавания образов.

Обобщенные линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая постановка задачи машинного обучения

Линейная регрессия

Задача классификации  
Логистическая регрессия

Метод IRLS

- Так как прямая минимизация правдоподобия невозможна, воспользуемся итерационным методом Ньютона
- Обоснованием корректности использования метода Ньютона является унимодальность оптимизируемой функции  $L(\mathbf{w})$  и ее гладкость во всем пространстве весов
- Формула пересчета в методе Ньютона

$$\mathbf{w}^{new} = \mathbf{w}^{old} - H^{-1} \nabla L(\mathbf{w}),$$

где  $H = \nabla \nabla L(\mathbf{w})$  — гессиан логарифма правдоподобия обучающей выборки

# Формулы пересчета

Обозначим  $s_i = \frac{1}{1 + \exp(-t_i y_i)}$ , тогда:

$$\nabla L(\mathbf{w}) = \Phi^T(\mathbf{s} - \mathbf{t}), \quad \nabla \nabla L(\mathbf{w}) = \Phi^T R \Phi$$

$$R = \begin{pmatrix} s_1(1-s_1) & 0 & \dots & 0 \\ 0 & s_2(1-s_2) & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & s_n(1-s_n) \end{pmatrix}$$

$$\mathbf{w}^{new} = \mathbf{w}^{old} - (\Phi^T R \Phi)^{-1} \Phi^T(\mathbf{s} - \mathbf{t}) =$$

$$(\Phi^T R \Phi)^{-1} (\Phi^T R \Phi \mathbf{w}^{old} - \Phi^T R R^{-1}(\mathbf{s} - \mathbf{t})) = (\Phi^T R \Phi)^{-1} \Phi^T R \mathbf{z},$$

где  $\mathbf{z} = \Phi \mathbf{w}^{old} - R^{-1}(\mathbf{s} - \mathbf{t})$

Название метода (метод наименьших квадратов с итеративно пересчитываемыми весами) связано с тем, что последняя формула является формулой для взвешенного МНК (веса задаются диагональной матрицей  $R$ ), причем на каждой итерации веса корректируются

# Заключительные замечания

Лекция 2.  
Вероятностная  
постановка  
задачи  
распознавания  
образов.  
Обобщенные  
линейные модели

Ветров,  
Кропотов

Ликбез

Статистическая  
постановка  
задачи  
машинного  
обучения

Линейная  
регрессия

Задача  
классификации  
Логистическая  
регрессия

Метод IRLS

- На практике матрица  $\Phi^T R \Phi$  часто бывает вырождена (всегда при  $m > n$ ), поэтому обычно прибегают к регуляризации матрицы  $(\Phi^T R \Phi + \lambda I)$
- !! Параметр регуляризации  $\lambda$  является структурным параметром!!
- !! Базисные функции  $\phi_j(\mathbf{x})$ , а значит и матрица  $\Phi$  являются структурными параметрами!!
- С поиском методов автоматического выбора базисных функций связана одна из наиболее интригующих проблем современного машинного обучения