

# РЕГУЛЯРИЗАЦИЯ В ЛИНЕЙНОЙ РЕГРЕССИИ

---

Сергей Николенко

uData School — Киев

1 июня 2018 г.

---

## *Random facts:*

- 1 июня 1838 г. на деньги, вырученные от продажи написанного Брюлловым портрета Жуковского, был выкуплен у помещика Энгельгардта крепостной Тарас Шевченко
- Международный день защиты детей 1 июня был учреждён в ноябре 1949 года в Париже решением конгресса Международной демократической федерации женщин под председательством физика Эжени Коттон, ученицы Марии Кюри, и впервые отмечался 1 июня 1950 года
- 1 июня 1965 г. Михаилу Шолохову была присуждена Нобелевская премия по литературе
- 1 июня 2001 г. наследный принц Непала Дипендра расстрелял всю свою семью, включая короля Бирендру, и застрелился сам

# РЕГРЕССИЯ ПО-БАЙЕСОВСКИ

---

- А теперь давайте посмотрим на регрессию с совсем байесовской стороны.
- Напомним основу байесовского подхода:
  1. найти апостериорное распределение на гипотезах/параметрах:

$$p(\theta | D) \propto p(D|\theta)p(\theta)$$

(и/или найти максимальную апостериорную гипотезу  $\arg \max_{\theta} p(\theta | D)$ );

2. найти апостериорное распределение исходов дальнейших экспериментов:

$$p(x | D) \propto \int_{\theta \in \Theta} p(x | \theta)p(D|\theta)p(\theta)d\theta.$$

- В нашем рассмотрении пока не было никаких априорных распределений.
- Давайте какое-нибудь введём; например, нормальное (почему так – позже):

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mu_0, \Sigma_0).$$

- Рассмотрим набор данных  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  со значениями  $\mathbf{t} = \{t_1, \dots, t_N\}$ . В этой модели мы предполагаем, что данные независимы и одинаково распределены:

$$p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^\top \phi(\mathbf{x}_n), \sigma^2).$$

- Тогда наша задача – посчитать

$$\begin{aligned} p(\mathbf{w} | \mathbf{t}) &\propto p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \sigma^2) p(\mathbf{w}) \\ &= \mathcal{N}(\mathbf{w} | \mu_0, \Sigma_0) \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^\top \phi(\mathbf{x}_n), \sigma^2). \end{aligned}$$

- Давайте подсчитаем.

- Получится

$$\begin{aligned}p(\mathbf{w} | \mathbf{t}) &= \mathcal{N}(\mathbf{w} | \mu_N, \Sigma_N), \\ \mu_N &= \Sigma_N \left( \Sigma_0^{-1} \mu_0 + \frac{1}{\sigma^2} \Phi^\top \mathbf{t} \right), \\ \Sigma_N &= \left( \Sigma_0^{-1} + \frac{1}{\sigma^2} \Phi^\top \Phi \right)^{-1}.\end{aligned}$$

- Теперь давайте подсчитаем логарифм правдоподобия.

- Если мы возьмём априорное распределение около нуля:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid 0, \frac{1}{\alpha} \mathbf{I}),$$

то логарифм правдоподобия получится

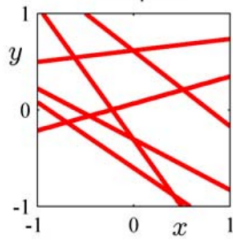
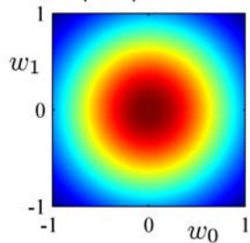
$$\ln p(\mathbf{w} \mid \mathbf{t}) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \phi(\mathbf{x}_n))^2 - \frac{\alpha}{2} \mathbf{w}^\top \mathbf{w} + \text{const},$$

то есть в точности гребневая регрессия.

likelihood

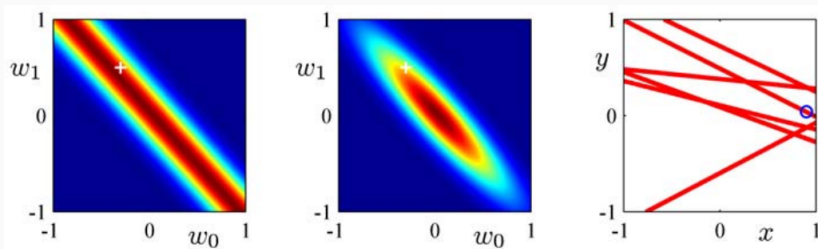
prior/posterior

data space

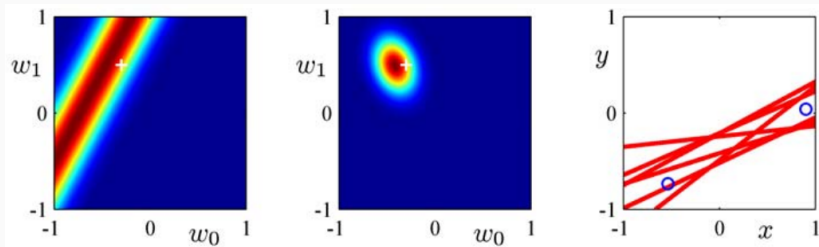




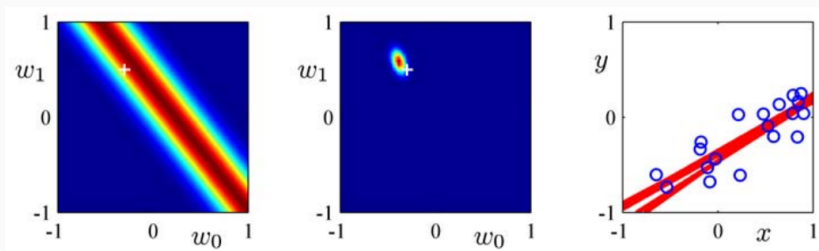
# ПРИМЕР



# ПРИМЕР



# ПРИМЕР



- Можно слегка обобщить – рассмотреть априорное распределение более общего вида

$$p(\mathbf{w} \mid \alpha) = \left[ \frac{q}{2} \left( \frac{\alpha}{2} \right)^{1/q} \frac{1}{\Gamma(1/q)} \right]^M e^{-\frac{\alpha}{2} \sum_{j=1}^M |w_j|^q}.$$

**Упражнение.** Подсчитайте логарифм правдоподобия.

## РАЗНЫЕ РЕГУЛЯРИЗАТОРЫ

---

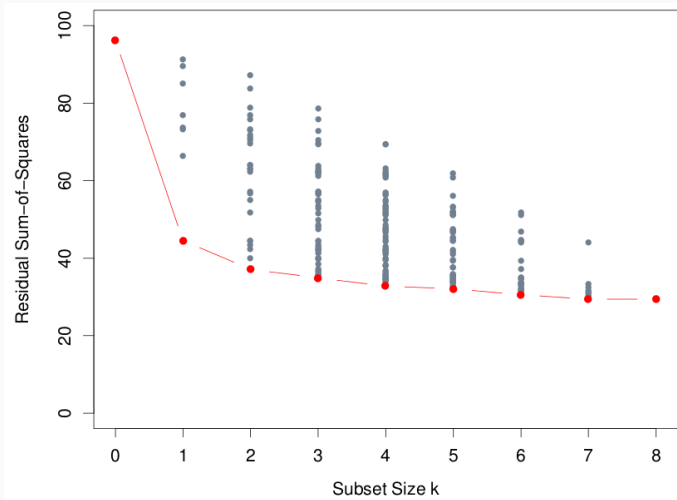
- Мы знаем, что наименьшие квадраты не всегда хорошо работают. Две причины:
  1. плохая предсказательная сила – часто лучше регуляризовать, пожертвовав  $\text{bias}$ 'ом в пользу  $\text{variance}$ ;
  2. сложности в интерпретации – хотелось бы понимать, что происходит, если переменных с ненулевыми коэффициентами слишком много, не получится.
- Мораль: хотелось бы сделать так, чтобы было поменьше ненулевых компонент в векторе  $\mathbf{w}$ .

- Может быть, давайте так и сделаем? Будем искать самые лучшие компоненты и делать их ненулевыми.
- Это называется subset selection.
- Можно просто делать best subset selection: выбирать подмножество из  $k$  входных переменных, которые дают самые лучшие результаты.

- Это долго, даже если делать с умом, потому что subsets много.
- Forward-stepwise selection: начинаем со свободного члена, потом добавляет на каждом шаге предиктор, который максимально уменьшает ошибку.
- Т.е. подмножества тут получаются вложенные.
- Backward-stepwise selection: начинаем с полной регрессии и на каждом шаге убираем предиктор, который оказывает меньше всего влияния на ошибку.



# SUBSET SELECTION



- Теперь давайте рассмотрим лассо-регрессию:

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (f(x_i, \mathbf{w}) - y_i)^2 + \lambda \sum_{j=0}^p |w_j|.$$

- Главное отличие – теперь форма ограничений (т.е. форма априорного распределения) такова, что весьма вероятно получить строго нулевые  $w_j$ .
- Кстати, что значит «форма ограничений»?

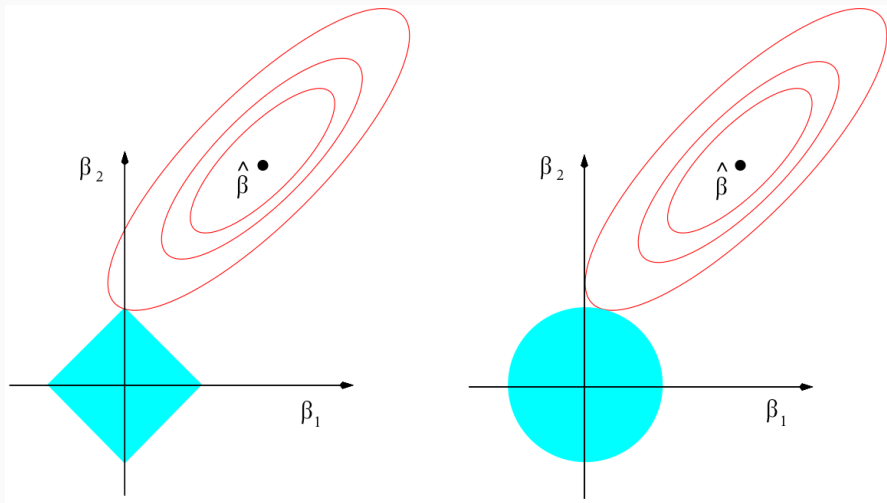
- Мы можем переписать регрессию с регуляризатором по-другому:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left\{ \frac{1}{2} \sum_{i=1}^N (f(x_i, \mathbf{w}) - y_i)^2 + \lambda \sum_{j=0}^p |w_j| \right\},$$

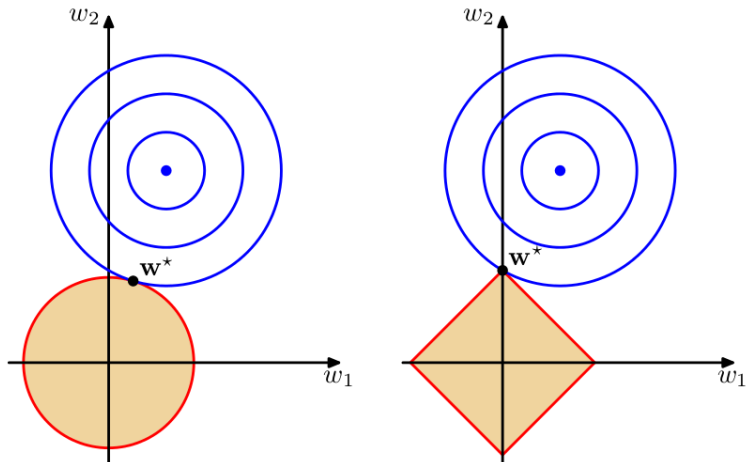
эквивалентно

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left\{ \frac{1}{2} \sum_{i=1}^N (f(x_i, \mathbf{w}) - y_i)^2 \right\} \text{ при } \sum_{j=0}^p |w_j| \leq t.$$

**Упражнение.** Докажите это.



# ГРЕБЕНЬ И ЛАССО

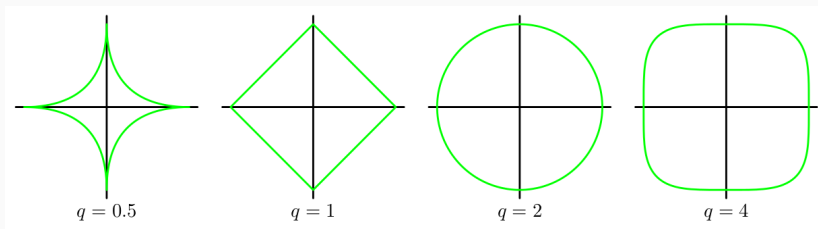
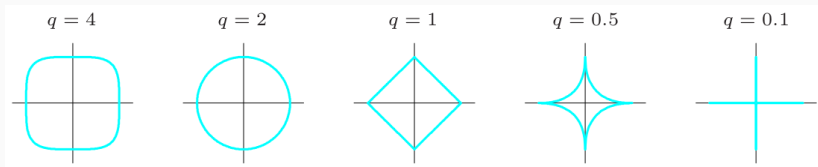


- Можно рассмотреть обобщение гребневой и лассо-регрессии:

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (f(x_i, \mathbf{w}) - y_i)^2 + \lambda \sum_{j=0}^p (|w_j|)^q.$$

**Упражнение.** Какому априорному распределению на параметры  $\mathbf{w}$  соответствует эта задача?

# РАЗНЫЕ $q$



## Упражнение.

1. Реализуйте линейную регрессию с одномерным входом и гауссовскими базисными функциями

$$\varphi_j(x) = e^{-\frac{1}{2s}(x-\mu_j)^2},$$

где  $\mu_j$  для  $j = 1..M$  распределены по отрезку равномерно.

2. Реализуйте разные виды регуляризации, нарисуйте графики полученных моделей для разных  $s$ ,  $M$  и  $\alpha$ .



# ПРЕДСКАЗАНИЯ В ЛИНЕЙНОЙ РЕГРЕССИИ

---

- Вспомним задачи байесовского вывода:
  1. найти апостериорное распределение на гипотезах/параметрах:

$$p(\theta | D) \propto p(D|\theta)p(\theta)$$

(и/или найти максимальную апостериорную гипотезу  $\arg \max_{\theta} p(\theta | D)$ );

2. найти апостериорное распределение исходов дальнейших экспериментов:

$$p(x | D) \propto \int_{\theta \in \Theta} p(x | \theta)p(D|\theta)p(\theta)d\theta.$$

- В прошлый раз мы нашли апостериорное распределение: для гауссовского априорного

$$p(\mathbf{w} \mid \alpha) = \mathcal{N}(\mathbf{w} \mid 0, \frac{1}{\alpha} \mathbf{I})$$

мы нашли

$$\begin{aligned} p(\mathbf{w} \mid \mathbf{t}, \alpha, \beta) &= \mathcal{N}(\mathbf{w} \mid \mu_N, \Sigma_N), \\ \mu_N &= \Sigma_N \left( \Sigma_0^{-1} \mu_0 + \beta \Phi^T \mathbf{t} \right), \\ \Sigma_N &= \left( \Sigma_0^{-1} + \beta \Phi^T \Phi \right)^{-1}, \end{aligned}$$

где  $\beta = \frac{1}{\sigma^2}$  (precision нормального распределения).

- Теперь сделаем следующий шаг – найдём апостериорное распределение наших предсказаний

$$p(t | \mathbf{t}, \alpha, \beta) = \int p(t | \mathbf{w}, \beta)p(\mathbf{w} | \mathbf{t}, \alpha, \beta)d\mathbf{w}.$$

- Это свёртка двух гауссианов, и получается...

- ...тоже гауссиан:

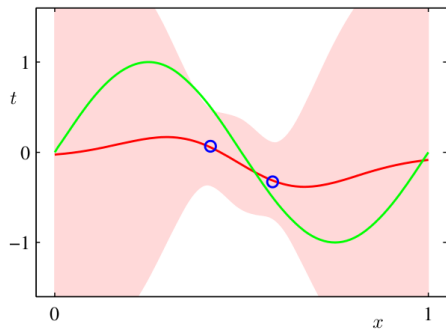
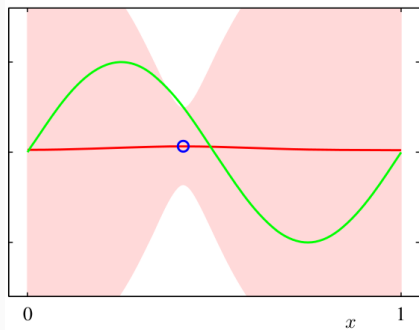
$$p(t \mid \mathbf{t}, \alpha, \beta) = \mathcal{N}(t \mid \mu_N^\top \phi(\mathbf{x}), \sigma_N^2),$$

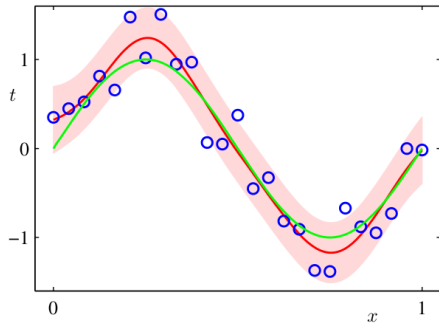
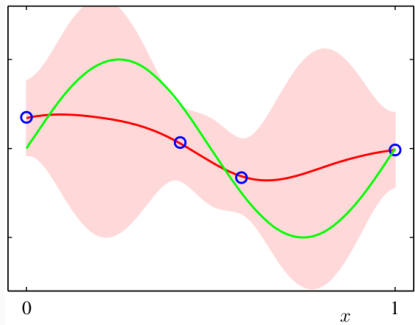
$$\text{где } \sigma_N^2 = \frac{1}{\beta} + \phi(\mathbf{x})^\top \Sigma_N \phi(\mathbf{x}).$$

- Т.е. дисперсия складывается из шума в данных  $\beta$  и дисперсии параметров  $\mathbf{w}$ ; гауссианы независимы, и их дисперсии просто складываются.

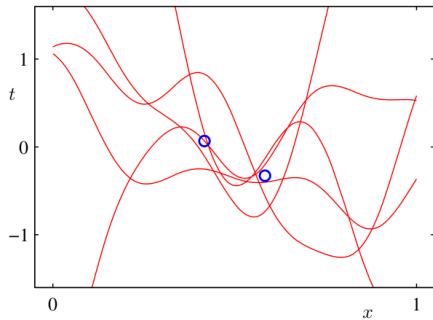
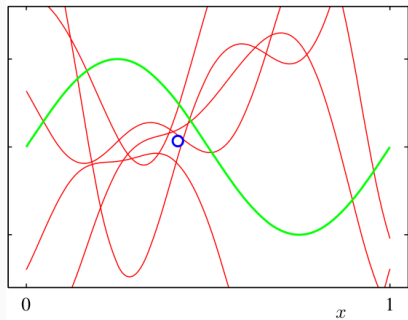
**Упражнение.** Оценка всё время уточняется:  $\sigma_{N+1}^2 \leq \sigma_N^2$ .

# ПРЕДСКАЗАНИЯ

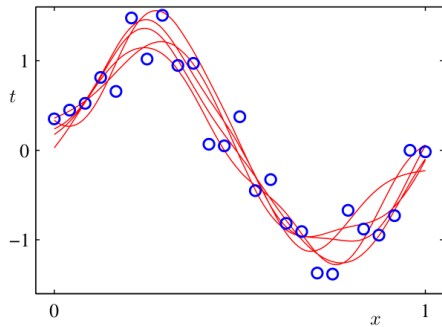
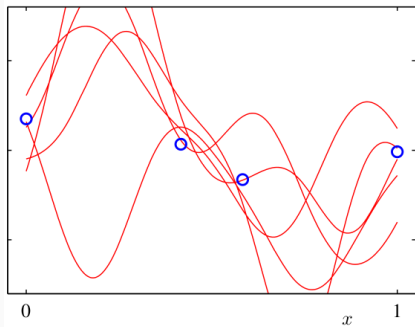




# ПРЕДСКАЗАНИЯ







Спасибо за внимание!