

Мат.модели машинного обучения: линейные модели регрессии и метод главных компонент

Воронцов Константин Вячеславович

`k.v.vorontsov@phystech.edu`

`http://www.MachineLearning.ru/wiki?title=User:Vokov`

Этот курс доступен на странице вики-ресурса

`http://www.MachineLearning.ru/wiki`

«Машинное обучение (курс лекций, К.В.Воронцов)»

1 Многомерная линейная регрессия

- Метод наименьших квадратов
- Многомерная линейная регрессия
- Сингулярное разложение

2 Регуляризация

- L_2 -регуляризация: гребневая регрессия
- L_1 -регуляризация: лассо Тибширани
- Негладкие регуляризаторы

3 Метод главных компонент

- Постановка задачи и основная теорема
- Метод главных компонент для линейной регрессии
- Обобщения

Метод наименьших квадратов (МНК)

- X — объекты (часто \mathbb{R}^n); Y — ответы (часто \mathbb{R} , реже \mathbb{R}^m);
 $X^\ell = (x_i, y_i)_{i=1}^\ell$ — обучающая выборка;
 $y_i = y(x_i)$, $y: X \rightarrow Y$ — неизвестная зависимость;
- $a(x, w)$ — модель зависимости,
 $w \in \mathbb{R}^p$ — вектор параметров модели.
- Метод наименьших квадратов (МНК):

$$Q(w, X^\ell) = \sum_{i=1}^{\ell} \gamma_i (a(x_i, w) - y_i)^2 \rightarrow \min_w,$$

где γ_i — вес, степень важности i -го объекта.

$Q(w^*, X^\ell)$ — остаточная сумма квадратов
 (residual sum of squares, RSS).

Многомерная линейная регрессия

$f_1(x), \dots, f_n(x)$ — числовые признаки;

Модель многомерной линейной регрессии:

$$a(x, w) = \sum_{j=1}^n w_j f_j(x), \quad w \in \mathbb{R}^n.$$

Матричные обозначения:

$$F_{\ell \times n} = \begin{pmatrix} f_1(x_1) & \dots & f_n(x_1) \\ \dots & \dots & \dots \\ f_1(x_\ell) & \dots & f_n(x_\ell) \end{pmatrix}, \quad y_{\ell \times 1} = \begin{pmatrix} y_1 \\ \dots \\ y_\ell \end{pmatrix}, \quad w_{n \times 1} = \begin{pmatrix} w_1 \\ \dots \\ w_n \end{pmatrix}.$$

Функционал квадрата ошибки:

$$Q(w, X^\ell) = \sum_{i=1}^{\ell} (a(x_i, w) - y_i)^2 = \|Fw - y\|^2 \rightarrow \min_w.$$

Нормальная система уравнений

Необходимое условие минимума в матричном виде:

$$\frac{\partial Q(w)}{\partial w} = 2F^T(Fw - y) = 0,$$

откуда следует *нормальная система* задачи МНК:

$$F^T F w = F^T y,$$

где $F^T F$ — матрица размера $n \times n$ полного ранга n .

Решение системы: $w^* = (F^T F)^{-1} F^T y = F^+ y$.

Значение функционала: $Q(w^*) = \|P_F y - y\|^2$,

где $P_F = F F^+ = F(F^T F)^{-1} F^T$ — *проекционная матрица*.

Геометрическая интерпретация МНК

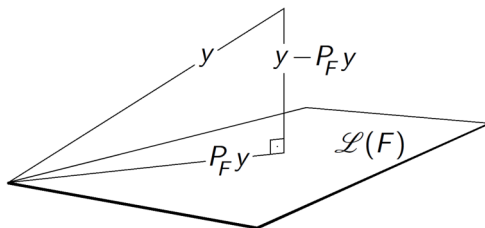
Линейная оболочка столбцов матрицы $F = (f_1, \dots, f_n)$, $f_j \in \mathbb{R}^\ell$:

$$\mathcal{L}(F) = \left\{ \sum_{j=1}^n w_j f_j \mid w \in \mathbb{R}^n \right\}$$

$P_F = F(F^T F)^{-1} F^T$ — проекционная матрица

$P_F y$ — проекция вектора $y \in \mathbb{R}^\ell$ на подпространство $\mathcal{L}(F)$

$(I_\ell - P_F)y$ — проекция y на его ортогональное дополнение



МНК — это опускание перпендикуляра в \mathbb{R}^ℓ из y на $\mathcal{L}(F)$

Сингулярное разложение

Произвольная $\ell \times n$ -матрица представима в виде *сингулярного разложения* (singular value decomposition, SVD):

$$F = VDU^T.$$

Основные свойства сингулярного разложения:

- 1 $\ell \times n$ -матрица $V = (v_1, \dots, v_n)$ ортогональна, $V^T V = I_n$, столбцы v_j — собственные векторы $\ell \times \ell$ -матрицы FF^T ;
- 2 $n \times n$ -матрица $U = (u_1, \dots, u_n)$ ортогональна, $U^T U = I_n$, столбцы u_j — собственные векторы $n \times n$ -матрицы $F^T F$;
- 3 $n \times n$ -матрица D диагональна, $D = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$, $\lambda_j \geq 0$ — общие собственные значения матриц $F^T F$ и FF^T .

Решение МНК через сингулярное разложение

Псевдообратная $F^+ = (F^T F)^{-1} F^T$, вектор МНК-решения w^* ,
 МНК-аппроксимация целевого вектора Fw^* :

$$F^+ = (UDV^T VDU^T)^{-1}UDV^T = UD^{-1}V^T = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j v_j^T;$$

$$w^* = F^+ y = UD^{-1}V^T y = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j (v_j^T y);$$

$$Fw^* = P_F y = (VDU^T)UD^{-1}V^T y = VV^T y = \sum_{j=1}^n v_j (v_j^T y);$$

$$\|w^*\|^2 = \|UD^{-1}V^T y\|^2 = \|D^{-1}V^T y\|^2 = \sum_{j=1}^n \frac{1}{\lambda_j} (v_j^T y)^2.$$

Тождества: $(AB)^{-1} = B^{-1}A^{-1}$, $(AB)^T = B^T A^T$, $\|w\|^2 = w^T w$

Мультиколлинеарность и число обусловленности матрицы

Если $\exists u \in \mathbb{R}^n$: $S_{n \times n} u \approx 0$, то некоторые с.з. S близки к нулю

Число обусловленности $n \times n$ -матрицы S :

$$\mu(S) = \|S\| \|S^{-1}\| = \frac{\max_{u: \|u\|=1} \|Su\|}{\min_{u: \|u\|=1} \|Su\|} = \frac{\lambda_{\max}}{\lambda_{\min}}$$

При умножении обратной матрицы на вектор, $z = S^{-1}u$, относительная погрешность усиливается в $\mu(S)$ раз:

$$\frac{\|\delta z\|}{\|z\|} \leq \mu(S) \frac{\|\delta u\|}{\|u\|}$$

В нашем случае: $S = F^T F$, $u = F^T y$, $w^* = S^{-1}u = (F^T F)^{-1} F^T y$, погрешности измерения признаков $f_j(x_i)$ и ответов y_i усиливаются в $\mu(F^T F)$ раз!

Стратегии устранения мультиколлинеарности

Если матрица $S = F^T F$ плохо обусловлена, то:

- решение w^* неустойчиво и плохо интерпретируемо, содержит большие по модулю w_j^* разных знаков;
- $\|w^*\|$ велико;
- возникает переобучение:
 на обучении $Q(w^*, X^\ell) = \|Fw^* - y\|^2$ мало;
 на контроле $Q(w^*, X^k) = \|F'w^* - y'\|^2$ велико;

Стратегии устранения мультиколлинеарности и переобучения:

- 1 регуляризация: $\|w\| \rightarrow \min$;
- 2 отбор признаков: $f_1, \dots, f_n \rightarrow f_{j_1}, \dots, f_{j_m}, \quad m \ll n$.
- 3 преобразование признаков: $f_1, \dots, f_n \rightarrow g_1, \dots, g_m, \quad m \ll n$;

Гребневая регрессия (ridge regression)

Штраф за увеличение L_2 -нормы вектора весов $\|w\|$:

$$Q_\tau(w) = \|Fw - y\|^2 + \frac{\tau}{2}\|w\|^2,$$

где τ — неотрицательный параметр регуляризации.

Модифицированное МНК-решение (τI_n — «гребень», ridge):

$$\frac{\partial Q_\tau(w)}{\partial w} = 2F^T(Fw - y) + 2\tau w = 0$$

$$w_\tau^* = (F^T F + \tau I_n)^{-1} F^T y.$$

Преимущество сингулярного разложения:

можно подбирать параметр τ , вычислив SVD только один раз.

Регуляризованный МНК через сингулярное разложение

Вектор регуляризованного МНК-решения w_τ^*
 и МНК-аппроксимация целевого вектора Fw_τ^* :

$$w_\tau^* = U(D^2 + \tau I_n)^{-1} D V^T y = \sum_{j=1}^n \frac{\sqrt{\lambda_j}}{\lambda_j + \tau} u_j (v_j^T y);$$

$$Fw_\tau^* = V D U^T w_\tau^* = V \operatorname{diag} \left(\frac{\lambda_j}{\lambda_j + \tau} \right) V^T y = \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau} v_j (v_j^T y);$$

$$\|w_\tau^*\|^2 = \|(D^2 + \tau I_n)^{-1} D V^T y\|^2 = \sum_{j=1}^n \frac{\lambda_j}{(\lambda_j + \tau)^2} (v_j^T y)^2.$$

$Fw_\tau^* \neq Fw^*$, но зато решение становится гораздо устойчивее.

Выбор параметра регуляризации τ

Контрольная выборка: $X^k = (x'_i, y'_i)_{i=1}^k$;

$$F'_{k \times n} = \begin{pmatrix} f_1(x'_1) & \dots & f_n(x'_1) \\ \dots & \dots & \dots \\ f_1(x'_k) & \dots & f_n(x'_k) \end{pmatrix}, \quad y'_{k \times 1} = \begin{pmatrix} y'_1 \\ \dots \\ y'_k \end{pmatrix}.$$

Вычисление функционала Q на контрольных данных T раз потребует $O(kn^2 + knT)$ операций:

$$Q(w_\tau^*, X^k) = \|F' w_\tau^* - y'\|^2 = \left\| \underbrace{F' U}_{k \times n} \operatorname{diag}\left(\frac{\sqrt{\lambda_j}}{\lambda_j + \tau}\right) \underbrace{V^T y}_{n \times 1} - y' \right\|^2.$$

Зависимость $Q(\tau)$ обычно имеет характерный минимум.

Регуляризация сокращает «эффективную размерность»

Сжатие (shrinkage) или сокращение весов (weight decay):

$$\|w_\tau^*\|^2 = \sum_{j=1}^n \frac{\lambda_j}{(\lambda_j + \tau)^2} (v_j^\top y)^2 < \|w^*\|^2 = \sum_{j=1}^n \frac{1}{\lambda_j} (v_j^\top y)^2.$$

Почему говорят о сокращении эффективной размерности?

Роль размерности играет след проекционной матрицы:

$$\text{tr } F(F^\top F)^{-1} F^\top = \text{tr } (F^\top F)^{-1} F^\top F = \text{tr } I_n = n.$$

При использовании регуляризации:

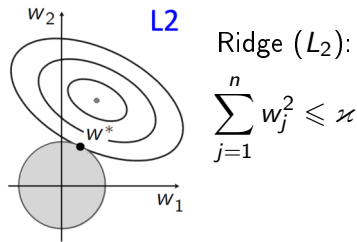
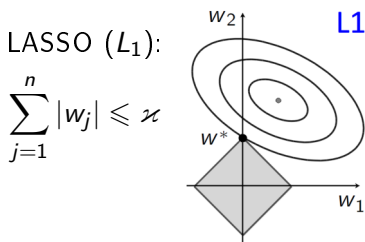
$$\text{tr } F(F^\top F + \tau I_n)^{-1} F^\top = \text{tr } \text{diag} \left(\frac{\lambda_j}{\lambda_j + \tau} \right) = \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau} < n.$$

Тождество: $\text{tr}(AB) = \text{tr}(BA)$

Регуляризация по L_1 -норме для отбора признаков

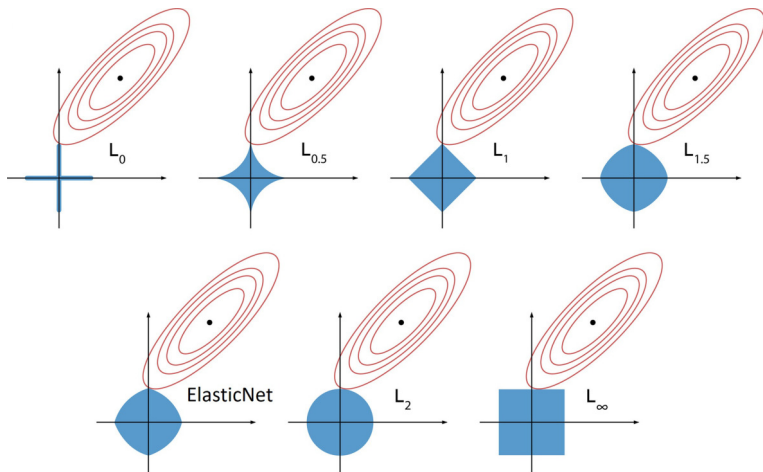
LASSO — Least Absolute Shrinkage and Selection Operator

$$\|Fw - y\|^2 + \mu \sum_{j=1}^n |w_j| \rightarrow \min_w \iff \begin{cases} \|Fw - y\|^2 \rightarrow \min_w; \\ \sum_{j=1}^n |w_j| \leq \kappa; \end{cases}$$



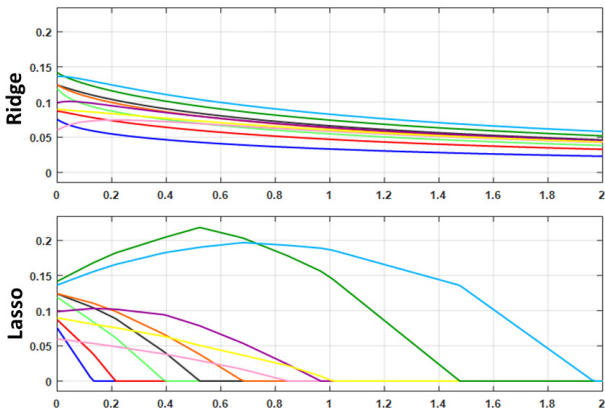
Геометрическая интерпретация отбора признаков

Сравнение регуляризаторов по различным L_p -нормам:



Сравнение L_2 (Ridge) и L_1 (LASSO) регуляризации

Типичный вид зависимости весов w_j от селективности μ



В LASSO с увеличением μ усиливается отбор признаков

Негладкие регуляризаторы для отбора и группировки признаков

Общий вид регуляризаторов (μ — параметр селективности):

$$\sum_{i=1}^{\ell} \mathcal{L}(w, x_i) + \sum_{j=1}^n R_{\mu}(w_j) \rightarrow \min_w .$$

Регуляризаторы с эффектами отбора и группировки признаков:

LASSO (L_1): $R_{\mu}(w) = \mu|w|$

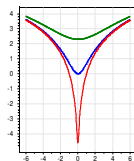
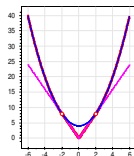
Elastic Net: $R_{\mu}(w) = \mu|w| + \tau w^2$

Support Feature Machine (SFM):

$$R_{\mu}(w) = \begin{cases} 2\mu|w|, & |w| \leq \mu; \\ \mu^2 + w^2, & |w| \geq \mu; \end{cases}$$

Relevance Feature Machine (RFM):

$$R_{\mu}(w) = \ln(\mu w^2 + 1)$$



Метод главных компонент (Principal Component Analysis, PCA)

$f_1(x), \dots, f_n(x)$ — исходные числовые признаки;

$g_1(x), \dots, g_m(x)$ — новые числовые признаки, $m \leq n$;

Требование: старые признаки $f_j(x)$ должны линейно восстанавливаться по новым признакам $g_s(x)$:

$$\hat{f}_j(x) = \sum_{s=1}^m g_s(x) u_{js}, \quad j = 1, \dots, n, \quad \forall x \in X,$$

как можно точнее на обучающей выборке x_1, \dots, x_ℓ :

$$\sum_{i=1}^{\ell} \sum_{j=1}^n (\hat{f}_j(x_i) - f_j(x_i))^2 \rightarrow \min_{\{g_s(x_i)\}, \{u_{js}\}}$$

Задача преобразования признаков (feature transformation)

— это **задача обучения без учителя**, тут нет ответов y_i

Матричные обозначения

Матрицы «объекты–признаки», старая и новая:

$$F_{\ell \times n} = \begin{pmatrix} f_1(x_1) & \dots & f_n(x_1) \\ \dots & \dots & \dots \\ f_1(x_\ell) & \dots & f_n(x_\ell) \end{pmatrix}; \quad G_{\ell \times m} = \begin{pmatrix} g_1(x_1) & \dots & g_m(x_1) \\ \dots & \dots & \dots \\ g_1(x_\ell) & \dots & g_m(x_\ell) \end{pmatrix}.$$

Матрица линейного преобразования новых признаков в старые:

$$U_{n \times m} = \begin{pmatrix} u_{11} & \dots & u_{1m} \\ \dots & \dots & \dots \\ u_{n1} & \dots & u_{nm} \end{pmatrix}; \quad \hat{F} = GU^T \stackrel{\text{ХОТИМ}}{\approx} F.$$

Найти: сразу и новые признаки G , и преобразование U :

$$\sum_{i=1}^{\ell} \sum_{j=1}^n (\hat{f}_j(x_i) - f_j(x_i))^2 = \|GU^T - F\|^2 \rightarrow \min_{G, U},$$

Основная теорема метода главных компонент

Теорема

Если $m \leq \text{rk } F$, то минимум $\|GU^T - F\|^2$ достигается, когда столбцы U — это с.в. матрицы $F^T F$, соответствующие m максимальным с.з. $\lambda_1, \dots, \lambda_m$, а матрица $G = FU$.

При этом:

- 1 матрица U ортонормирована: $U^T U = I_m$;
- 2 матрица G ортогональна: $G^T G = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$;
- 3 $U\Lambda = F^T F U$; $G\Lambda = FF^T G$;
- 4 $\|GU^T - F\|^2 = \|F\|^2 - \text{tr } \Lambda = \sum_{j=m+1}^n \lambda_j$.

Тождество: $\|F\|^2 = \text{tr}(F^T F) = \sum_{j=1}^n \lambda_j$

Связь с сингулярным разложением

Произвольная $\ell \times n$ -матрица F представима в виде SVD:

$$F = VDU^T; \quad U^T U = I_m; \quad V^T V = I_m; \quad D = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$$

Если взять $m = n$, то:

① $\|GU^T - F\|^2 = 0$

② представление $\hat{F} = GU^T = F$ точное и совпадает с SVD, если положить $G = V\sqrt{\Lambda}$, $D = \sqrt{\Lambda}$:

$$F = GU^T = V\sqrt{\Lambda}U^T;$$

③ линейное преобразование U работает в обе стороны:

$$F = GU^T; \quad G = FU.$$

$G^T G = \Lambda$ — новые признаки некоррелированы,

U — декоррелирующее преобразование Карунена–Лоэва

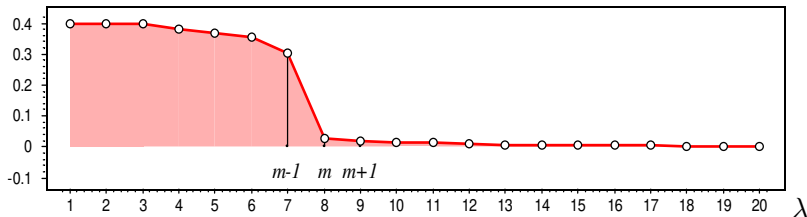
Эффективная размерность выборки

Упорядочим с.з. $F^T F$ по убыванию: $\lambda_1 \geq \dots \geq \lambda_n \geq 0$.

Эффективная размерность выборки — это наименьшее целое m , при котором

$$E_m = \frac{\|GU^T - F\|^2}{\|F\|^2} = \frac{\lambda_{m+1} + \dots + \lambda_n}{\lambda_1 + \dots + \lambda_n} \leq \varepsilon.$$

Критерий «крутого склона»: находим m : $E_{m-1} \gg E_m$:



Решение задачи НК в новых признаках

Заменим $F_{\ell \cdot n}$ на её приближение $G \cdot U^T$, предполагая $m \leq n$:

$$\|G \underbrace{U^T w}_z - y\|^2 = \|Gz - y\|^2 \rightarrow \min_z.$$

Связь нового и старого вектора коэффициентов:

$$z = U^T w; \quad w = Uz.$$

Решение задачи наименьших квадратов относительно z :

$$z^* = D^{-1} V^T y; \quad w^* = UD^{-1} V^T y = \sum_{j=1}^m \frac{1}{\sqrt{\lambda_j}} u_j (v_j^T y);$$

$$Gz^* = VV^T y = \sum_{j=1}^m v_j (v_j^T y);$$

Единственное отличие от прежнего w^* — m слагаемых вместо n

Спектральный метод наименьших квадратов

1. Построить SVD-разложение, упорядочить $\lambda_1 \geq \dots \geq \lambda_n$
2. Отделить $n - m$ наименьших с.з. от нуля: $\lambda'_j := \lambda_j + \delta_j$

Частные случаи:

- $\lambda'_j := \lambda_j + \tau$ — гребневая регрессия
- $\lambda'_j := \lambda_j + \infty [j > m]$ — метод главных компонент
- $\lambda'_j := \lambda_j + \tau [j > m]$ — нечто промежуточное

3. Применить формулы SVD для модификации МНК-решения:

$$w^* = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j (v_j^T y) \quad \longrightarrow \quad w^* = \sum_{j=1}^n \frac{\sqrt{\lambda_j}}{\lambda'_j} u_j (v_j^T y)$$

$$Fw^* = \sum_{j=1}^n v_j (v_j^T y) \quad \longrightarrow \quad Fw^* = \sum_{j=1}^n \frac{\lambda_j}{\lambda'_j} v_j (v_j^T y)$$

Задачи низкорангового матричного разложения

- Генерация новых признаков $f_1, \dots, f_n \rightarrow g_1, \dots, g_m$, $m \ll n$;
- Снижение размерности вектора признаков
- Восстановление пропущенных значений (missing values)

Дано: матрица $F = (f_{ij})_{\ell \times n}$, $(i, j) \in \Omega \subseteq \{1..l\} \times \{1..n\}$

Найти: матрицы $G = (g_{is})_{\ell \times m}$ и $U^T = (u_{sj})_{m \times n}$

Критерий: $\|F - GU^T\|^2 = \sum_{(i,j) \in \Omega} \left(f_{ij} - \sum_s g_{is} u_{sj} \right)^2 \rightarrow \min_{X, Y}$

Когда вместо SVD применяется SG или другие методы:

- разреженные данные: $|\Omega| \ll \ell n$
- неквадратичная функция потерь
- неотрицательное матричное разложение: $g_{is} \geq 0$, $u_{sj} \geq 0$

Примеры задач неотрицательного матричного разложения

- 1 **Выявление интересов в рекомендательных системах (recommender systems, collaborative filtering)**

$$f_{iu} = \sum_t p_{it} q_{tu}$$

дано: f_{iu} — рейтинги товаров i , поставленные пользователем u

найти: p_{it} — профиль интересов товара i

q_{tu} — профиль интересов пользователя u

- 2 **Разделение смеси химических веществ по данным жидкостной хроматографии**

$$f_{t\lambda} = \sum_i c_{ti} s_{i\lambda}$$

дано: $f_{t\lambda}$ — концентрация на выходе УФ-детектора

найти: c_{ti} — хроматограмма i -го вещества, t — время

$s_{i\lambda}$ — спектр i -го вещества, λ — длина волны

Примеры задач неотрицательного матричного разложения

- 3 Латентный семантический анализ коллекций текстов (тематическое моделирование)

$$f_{wd} = \sum_t \varphi_{wt} \theta_{td}$$

дано: $f_{wd} = p(w|d)$ — частоты слов w в документах d

найти: $\varphi_{wt} = p(w|t)$ — распределения слов w в темах t

$\theta_{td} = p(t|d)$ — распределения тем t в документах d

- 4 Оценивание экспрессии генов по данным ДНК-микрочипов с учётом кросс-гибридизации

$$f_{pk} = \sum_g a_{pg} c_{gk}$$

дано: f_{pk} — интенсивность свечения p -й пробы на k -м чипе

найти: a_{pg} — коэффициент сродства p -й пробы g -му гену

c_{gk} — концентрация g -го гена на k -м чипе

- Многомерная линейная регрессия
— через *сингулярное разложение*
- Три приёма против мультиколлинеарности и переобучения
— регуляризация, отбор и преобразование признаков
- L_2 -регуляризация, она же гребневая регрессия
— тоже через *сингулярное разложение*
- L_1 -регуляризация (LASSO) и др. негладкие регуляризаторы
— регулируемый отбор признаков
- Метод главных компонент — задача матричного разложения
— снова через *сингулярное разложение*
- Другие методы матричных разложений и их приложения
— в следующем семестре