$$\min_{x \in X \subseteq \mathbb{R}^d} f(x)$$

$X$ — „простое" мн-во

# Example: Master's Admission

- $0.0 \leq$ GPA $\leq 4.0$ (from F to A)
- $0 \leq$ Salary
- $1.0 \leq$ Perfomance $\leq 6.0$ (final score of secess)
- Historical data:

| GPA | Salary | Perfomance |
|---|---|---|
| 3.52 | 100 | 3.92 |
| 3.66 | 109 | 4.34 |
| 3.76 | 113 | 4.80 |
| 3.74 | 100 | 4.67 |
| 3.93 | 100 | 5.52 |
| 3.88 | 115 | 5.44 |
| 3.77 | 115 | 5.04 |
| 3.66 | 107 | 4.73 |
| 3.87 | 106 | 5.03 |
| 3.84 | 107 | 5.06 |

# Master's Admission: Linear model

Hypothesis:

$$\text{Perfomance} \approx w_0 + w_1 \cdot \text{GPA} + w_2 \cdot \text{Salary}$$

for weights $w_0, w_1, w_2$ to be learned.

Approach: Find $w_0, w_1, w_2$ by minimizing least squares error over the historical data.

**Question:** what we need to do with data before solving something?

# Master's Admission: Linear model

Hypothesis:

$$\text{Perfomance} \approx w_0 + w_1 \cdot \text{GPA} + w_2 \cdot \text{Salary}$$

for weights $w_0, w_1, w_2$ to be learned.

Approach: Find $w_0, w_1, w_2$ by minimizing least squares error over the historical data.

**Question:** what we need to do with data before solving something?

- Relevant GPA scores span a range of 0.5 (take only top students).
- Relevant Salary scores span a range of 20 (from 100 to 120 - others go to jobs, not to master).
- $\Rightarrow$ normalize first so that $w_1, w_2$ can be compared

# General setting

$n$ inputs $x_1, \ldots, x_n$, $x_i \in \mathbb{R}^d$ for all $i$

$d$ input variables $1, 2, \ldots, d$

- 10 (GPA, Salary) pairs, two input variables

$n$ outputs $y_1, \ldots, y_n \in \mathbb{R}$

- 10 Perfomance scores

$(x_i, y_i)$: an observation

- $((3.93, 100), 5.52)$, observation (of a student doing very well)

With weights $w_0, w = (w_1, \ldots, w_d) \in \mathbb{R}^d$, we plan to minimize the least squares objective

$$f(w_0, w) = \sum_{i=1}^{n} (w_0 + w^T x_i - y_i)^2.$$

# General setting: centering

Want to assume that

$$\frac{1}{n}\sum_{i=1}^{n}\mathsf{x}_i = 0, \quad \frac{1}{n}\sum_{i=1}^{n}y_i = 0.$$

Can be achieved by

- subtracting the mean $\bar{\mathsf{x}} = \frac{1}{n}\sum_{i=1}^{n}\mathsf{x}_i$ from every input
- subtracting the mean $\bar{y} = \frac{1}{n}\sum_{i=1}^{n}y_i$ from every output.

**Question:** after centering what we can assume?

# General setting: centering

Want to assume that

$$\frac{1}{n}\sum_{i=1}^{n}x_i = 0, \quad \frac{1}{n}\sum_{i=1}^{n}y_i = 0.$$

Can be achieved by

- subtracting the mean $\bar{x} = \frac{1}{n}\sum_{i=1}^{n}x_i$ from every input
- subtracting the mean $\bar{y} = \frac{1}{n}\sum_{i=1}^{n}y_i$ from every output.

**Question:** after centering what we can assume?

After centering: $w_0^* = 0, w^*$ is unaffected

$\Rightarrow$ From now on consider function

$$f(w) = \sum_{i=1}^{n}(w^T x_i - y_i)^2.$$

# General setting: normalization

Want to assume that for all $j$, the $n$ input values $x_{1j}, \ldots x_{nj}$ are on the same scale:

$$\frac{1}{n} \sum_{i=1}^{n} x_{ij}^2 = 1, \quad j = 1, \ldots, d.$$

Can be achieved by

- multiplying $x_{ij}$ by $s(j) = \sqrt{n / \sum_{i=1}^{n} x_{ij}^2}$ for all $i, j$
- in $w^*$, this just multiplies $w_j^*$ by $1/s(j)$

**Constrained Optimization**
ooooooooo●oooo

Projection Gradient Descent
ooooooooo

Frank-Wolfe Method
ooooooo

# Master's Admission: Centered and normalized data

| $x_{i1}$ (GPA) | $x_{i2}$ (Salary) | $y_i$ (Perfomance) |
|---:|---:|---:|
| -2.04 | -1.28 | -0.94 |
| -0.88 | 0.32 | -0.52 |
| -0.05 | 1.03 | -0.05 |
| -0.16 | -1.28 | -0.18 |
| 1.42 | -1.28 | 0.67 |
| 1.02 | 1.39 | 0.59 |
| 0.06 | 1.39 | 0.19 |
| -0.88 | -0.04 | -0.12 |
| 0.89 | -0.21 | 0.17 |
| 0.62 | -0.04 | 0.21 |

Least-squares objective:

$$\min \quad f(w_1, w_2) = \sum_{i=1}^{10} (w_1 x_{i1} + w_2 x_{i2} - y_i)^2.$$

# Master's Admission: Results

Optimal solution: $\min$

$(w_1, w_2) \in \mathbb{R}^2$

$$w^* = (w_1^*, w_2^*) \approx (0.43, 0.097)$$

# Master's Admission: Results

Optimal solution: $w^* = (w_1^*, w_2^*) \approx (0.43, 0.097)$
Under hypothesis (linear model), we expect $y_i \approx y_i^* = 0.43 x_{i1} + 0.097 x_{i2}$

| $x_{i1}$ | $x_{i2}$ | $y_i$ | $y_i^*$ |
|---------:|---------:|------:|--------:|
| -2.04 | -1.28 | -0.94 | -1.00 |
| -0.88 | 0.32 | -0.52 | -0.35 |
| -0.05 | 1.03 | -0.05 | 0.08 |
| -0.16 | -1.28 | -0.18 | -0.19 |
| 1.42 | -1.28 | 0.67 | 0.49 |
| 1.02 | 1.39 | 0.59 | 0.57 |
| 0.06 | 1.39 | 0.19 | 0.16 |
| -0.88 | -0.04 | -0.12 | -0.38 |
| 0.62 | -0.04 | 0.21 | 0.26 |

**Questiob:** what we can say about results? Salary has only very small influence ($w_2^* = 0.097$)

# Predicting Perfomance in the future

Problems:

- least squares solution is optimized for the training data, not for the future ("overfitting")
- "unimportant" variables should have weight 0, but they typically don't

**Constrained Optimization**
oooooooooo●oo

Projection Gradient Descent
ooooooooo

Frank-Wolfe Method
ooooooo

# Predicting Perfomance in the future

Problems:

- least squares solution is optimized for the training data, not for the future ("overfitting")
- "unimportant" variables should have weight 0, but they typically don't

**Subset selection heuristics**: drop variables with seemingly "small" contribution

**Constrained Optimization**
ooooooooo●oo

Projection Gradient Descent
ooooooooo

Frank-Wolfe Method
ooooooo

# Predicting Perfomance in the future

Problems:

- least squares solution is optimized for the training data, not for the future ("overfitting")

- "unimportant" variables should have weight 0, but they typically don't

**Subset selection heuristics**: drop variables with seemingly "small" contribution (various methods to decide what "small" means, and how many to drop)

**Best subset selection:** solve least squares subject to an additional constraint that there are at most $k$ nonzero weights. **Easy of not?**

$$n \qquad k \qquad C_n^k$$

**Constrained Optimization**
○○○○○○○○○●○○

Projection Gradient Descent
○○○○○○○○

Frank-Wolfe Method
○○○○○○○

# Predicting Perfomance in the future

Problems:

- least squares solution is optimized for the training data, not for the future ("overfitting")

- "unimportant" variables should have weight 0, but they typically don't

**Subset selection heuristics**: drop variables with seemingly "small" contribution (various methods to decide what "small" means, and how many to drop)

**Best subset selection**: solve least squares subject to an additional constraint that there are at most $k$ nonzero weights. **Easy of not?** Non-convex or NP-hard – various $k$ might have to be tried.

**Question:** if we have 100 features, how many different subsets (of features) can we have?

**Constrained Optimization**
ooooooooo●oo

Projection Gradient Descent
ooooooooo

Frank-Wolfe Method
ooooooo

# Predicting Perfomance in the future

Problems:

- least squares solution is optimized for the training data, not for the future ("overfitting")

- "unimportant" variables should have weight 0, but they typically don't

**Subset selection heuristics**: drop variables with seemingly "small" contribution (various methods to decide what "small" means, and how many to drop)

**Best subset selection**: solve least squares subject to an additional constraint that there are at most $k$ nonzero weights. **Easy of not?** Non-convex or NP-hard – various $k$ might have to be tried.

**Question:** if we have 100 features, how many different subsets (of features) can we have? $2^{100} \approx 1.26 \cdot 10^{30}$.

**LASSO**: popular approach with some favorable statistical properties

**Constrained Optimization**
oooooooooooo●o

Projection Gradient Descent
ooooooooo

Frank-Wolfe Method
ooooooo

## The LASSO: a constrained optimization problem

$$
\begin{array}{ll}
\text{minimize} & \sum_{i=1}^{n} \|w^\top x_i - y_i\|^2 \\
\text{subject to} & \|w\|_1 \leq R,
\end{array}
\tag{1}
$$

where $R \in \mathbb{R}_+$ is some parameter.

**Constrained Optimization**
ooooooooooooo●o

Projection Gradient Descent
ooooooooo

Frank-Wolfe Method
ooooooo

# The LASSO: a constrained optimization problem

$$
\begin{aligned}
\text{minimize} \quad & \sum_{i=1}^{n} \|w^\top x_i - y_i\|^2 \\
\text{subject to} \quad & \|w\|_1 \leq R,
\end{aligned}
\tag{1}
$$

where $R \in \mathbb{R}_+$ is some parameter.
$\|w\|_1 = \sum_{i=1}^{d} |w_j|$ is the 1-norm.

**Constrained Optimization**
ooooooooooo●o

Projection Gradient Descent
ooooooooo

Frank-Wolfe Method
ooooooo

## The LASSO: a constrained optimization problem

$$\begin{array}{ll} \text{minimize} & \sum_{i=1}^{n} \|\mathbf{w}^{\top}\mathbf{x}_i - y_i\|^2 \\ \text{subject to} & \|\mathbf{w}\|_1 \leq R, \end{array} \qquad (1)$$

where $R \in \mathbb{R}_+$ is some parameter.

$\|\mathbf{w}\|_1 = \sum_{i=1}^{d} |w_j|$ is the 1-norm.

In our case:

$R = 0.2 \Rightarrow \mathbf{w}^* = (w_1^*, w_2^*) = (0.2, 0)$:

**Constrained Optimization**
ooooooooooo●o

Projection Gradient Descent
ooooooooo

Frank-Wolfe Method
ooooooo

## The LASSO: a constrained optimization problem

$$\begin{array}{ll} \text{minimize} & \sum_{i=1}^{n} \|w^{\top} x_i - y_i\|^2 \\ \text{subject to} & \|w\|_1 \leq R, \end{array} \tag{1}$$

where $R \in \mathbb{R}_+$ is some parameter.

$\|w\|_1 = \sum_{i=1}^{d} |w_j|$ is the 1-norm.

In our case:

$R = 0.2 \Rightarrow \underline{w^* = (w_1^*, w_2^*) = (0.2, 0)}$: Salary is gone!

**Constrained Optimization**
○○○○○○○○○○○●○

Projection Gradient Descent
○○○○○○○○

Frank-Wolfe Method
○○○○○○○

# The LASSO: a constrained optimization problem

$$
\begin{aligned}
\text{minimize} \quad & \sum_{i=1}^{n} \|w^\top x_i - y_i\|^2 \\
\text{subject to} \quad & \|w\|_1 \leq R,
\end{aligned}
\tag{1}
$$

where $R \in \mathbb{R}_+$ is some parameter.

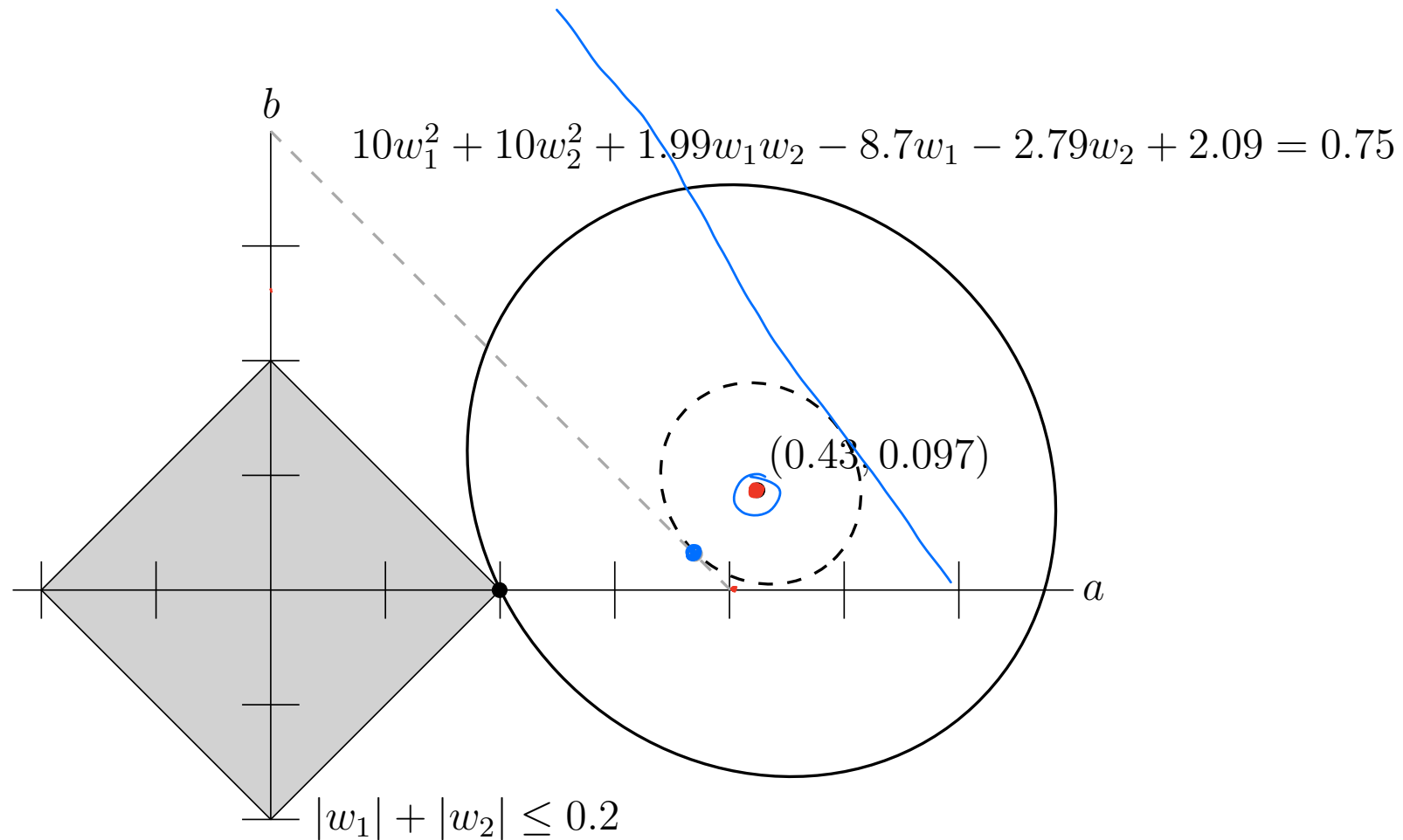$\|w\|_1 = \sum_{i=1}^{d} |w_j|$ is the 1-norm.

In our case:

$R = 0.2 \Rightarrow w^* = (w_1^*, w_2^*) = (0.2, 0)$: Salary is gone!

$R = 0.3 \Rightarrow w^* = (w_1^*, w_2^*) = (0.3, 0)$

**Constrained Optimization**
○○○○○○○○○○●○

Projection Gradient Descent
○○○○○○○○

Frank-Wolfe Method
○○○○○○○

# The LASSO: a constrained optimization problem

$$\begin{array}{ll} \text{minimize} & \sum_{i=1}^{n} \|w^{\top}x_i - y_i\|^2 \\ \text{subject to} & \|w\|_1 \le R, \end{array} \tag{1}$$

where $R \in \mathbb{R}_+$ is some parameter.

$\|w\|_1 = \sum_{i=1}^{d} |w_j|$ is the 1-norm.

In our case:

$R = 0.2 \Rightarrow w^* = (w_1^*, w_2^*) = (0.2, 0)$: Salary is gone!

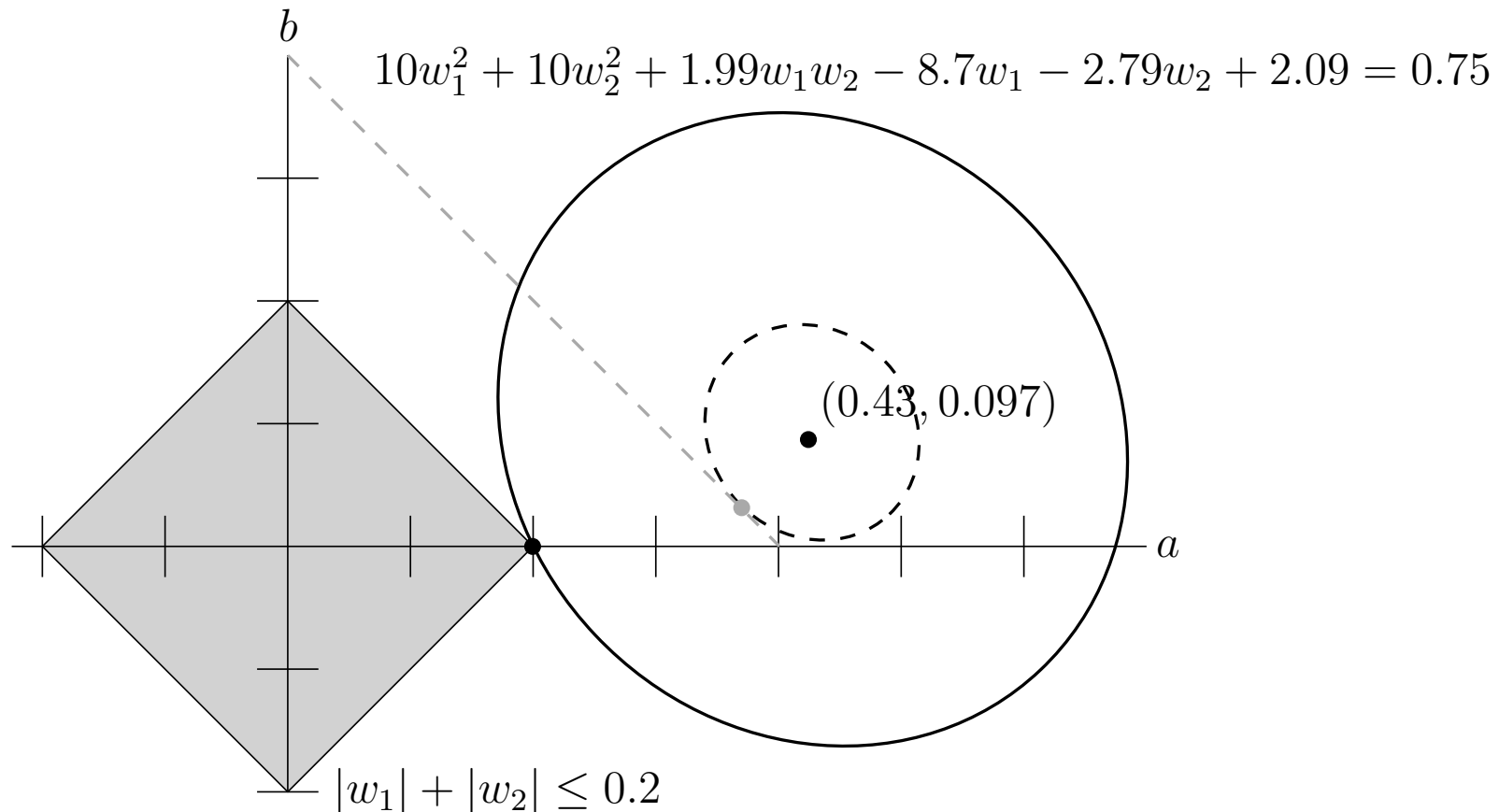$R = 0.3 \Rightarrow w^* = (w_1^*, w_2^*) = (0.3, 0)$

$R = 0.4 \Rightarrow w^* = (w_1^*, w_2^*) = (0.36, 0.036)$

**Constrained Optimization**
○○○○○○○○○○○●○

Projection Gradient Descent
○○○○○○○○

Frank-Wolfe Method
○○○○○○○

# The LASSO: a constrained optimization problem

$$\begin{aligned} \text{minimize} \quad & \sum_{i=1}^{n} \|w^{\top} x_i - y_i\|^2 \\ \text{subject to} \quad & \|w\|_1 \leq R, \end{aligned} \tag{1}$$

where $R \in \mathbb{R}_+$ is some parameter.

$\|w\|_1 = \sum_{i=1}^{d} |w_j|$ is the 1-norm.

In our case:

$R = 0.2 \Rightarrow w^* = (w_1^*, w_2^*) = (0.2, 0)$: Salary is gone!

$R = 0.3 \Rightarrow w^* = (w_1^*, w_2^*) = (0.3, 0)$

$R = 0.4 \Rightarrow w^* = (w_1^*, w_2^*) = (0.36, 0.036)$

$R \geq 0.6 \Rightarrow w^* = (w_1^*, w_2^*) = (0.43, 0.097)$

# Geometry of the LASSO



$$10w_1^2 + 10w_2^2 + 1.99w_1w_2 - 8.7w_1 - 2.79w_2 + 2.09 = 0.75$$

$(0.43, 0.097)$

$|w_1| + |w_2| \leq 0.2$

# Geometry of the LASSO



$$10w_1^2 + 10w_2^2 + 1.99w_1w_2 - 8.7w_1 - 2.79w_2 + 2.09 = 0.75$$

$(0.43, 0.097)$

$|w_1| + |w_2| \leq 0.2$

**Question:** Can we somehow modify gradient method to work with constraints?

## Условие оптимальности

- $f$ - непрер. диффер. на $\mathbb{R}^d$     $f: \mathbb{R}^d \to \mathbb{R}^d$
- $f$ - выпуклая функция
- $X$ - выпуклое

$x^* \in X$ - глобальный мин. $\min\limits_{x \in X} f(x) \iff$

$$\langle \nabla f(x^*); x - x^* \rangle \geq 0 \qquad \forall x \in X$$

### Физ. смысл:



$x^*$ - решение $\min\limits_{X} f(x)$

угол между $\nabla f(x^*); x - x^*$ - острый

## Док-во:

- достаточное. $\Longleftarrow$

$\langle \nabla f(x^*); x - x^* \rangle \geq 0 \qquad \forall x \in X$

выпуклость $f$:

$$f(x) \geq f(x^*) + \underbrace{\langle \nabla f(x^*); x - x^* \rangle}_{\geq 0} \geq f(x^*) \qquad \forall x \in X$$

$x^*$ - глоб. минимум на $X$

- необходимость. $\Longrightarrow$

$x^*$ - глобальный минимум на $X$

От противного: $\exists x \in X: \quad \langle \nabla f(x^*); x - x^* \rangle < 0$

$$x_\lambda = \lambda x + (1 - \lambda) x^* \qquad \lambda \in [0; 1]$$

$$\phi(\lambda) = f(x_\lambda) = f(\lambda x + (1 - \lambda) x^*)$$

$$\frac{d\phi}{d\lambda} = \frac{d}{d\lambda}\left(f(\lambda(x-x^*)+x^*)\right) = \langle \nabla f(\lambda(x-x^*)+x^*); x-x^* \rangle$$

$$\left.\frac{d\phi}{d\lambda}\right|_{\lambda=0} = \langle \nabla f(x^*); x-x^* \rangle \underbrace{< 0}_{\text{по пред.}}$$

$\phi$ убывает в окр. $0$, а значит $\exists \lambda > 0$ :

$$f(x^* + \lambda(x-x^*)) = \phi(\lambda) < \phi(0) = f(x^*)$$

противоречие $x^*$ — глоб. миним.

---

## Метод град. спуска с проекцией

$$x^{k+1} = x^k - \gamma \nabla f(x^k)$$

$\in X$

не сразу $x^{k+1} \in X$

$$\boxed{\begin{array}{l} x^{k+1} = \Pi_X\left(x^k - \gamma \nabla f(x^k)\right) \\[2mm] \Pi_X(y) = \underset{x \in X}{\arg\min} \|x-y\|_2^2 \quad \leftarrow \text{проекция (евклидова)} \end{array}}$$

## Св-ва проекции:

- $X$ — выпуклое, $x \in X$, $y \in \mathbb{R}^d$, тогда

$$\langle x - \Pi_X(y), y - \Pi_X(y) \rangle \le 0$$

Док-во: $\Pi_X(y) = \underset{z \in X}{\arg\min}\ d(z) := \underbrace{\|z-y\|_2^2}_{\text{выпуклая}}$

выпуклое

Условие оптимальности для $d(z)$ и $\Pi_X(y)$:

$$\langle \nabla d(z^*); z - z^* \rangle \ge 0 \qquad \forall z \in X$$

$$\langle \nabla d(\Pi_X(y)); \ x - \Pi_X(y) \rangle \geq 0$$

$$\nabla d(z) = 2(z-y)$$

$$2\langle \Pi_X(y) - y; \ x - \Pi_X(y) \rangle \geq 0 \quad \blacksquare$$

⊙ Нерастяжаемость оператора проекции

$X$ — выпуклое, $X_1, X_2 \in \mathbb{R}^d$, тогда

$$\|\Pi_X(x_1) - \Pi_X(x_2)\|_2 \leq \|x_1 - x_2\|_2$$

Док-во: пред. св-во $y = x_1, \ x = \Pi_X(x_2)$

$$\langle \Pi_X(x_2) - \Pi_X(x_1), x_1 - \Pi_X(x_1) \rangle \leq 0$$

аналогично $y = x_2, \ x = \Pi_X(x_1)$

$$\langle \Pi_X(x_1) - \Pi_X(x_2), x_2 - \Pi_X(x_2) \rangle \leq 0$$

складываем

$$\langle \Pi_X(x_2) - \Pi_X(x_1); \ x_1 - \Pi_X(x_1) - x_2 + \Pi_X(x_2) \rangle \leq 0$$

$$\langle \Pi_X(x_2) - \Pi_X(x_1); \ \Pi_X(x_2) - \Pi_X(x_1) \rangle \leq \langle x_2 - x_1; \ \Pi(x_2) - \Pi(x_1) \rangle$$

$$\|\Pi(x_2) - \Pi(x_1)\|_2^2 \leq \langle x_2 - x_1; \ \Pi(x_2) - \Pi(x_1) \rangle$$

КБШ

$$\|\Pi(x_2) - \Pi(x_1)\|_2^2 \leq \|\Pi(x_2) - \Pi(x_1)\|_2 \ \|x_2 - x_1\|_2$$

$$\|\Pi(x_2) - \Pi(x_1)\|_2 \leq \|x_1 - x_2\|_2 \quad \blacksquare$$

⊙ Стац. точки град. спуска с проекцией

$$x^* = \Pi_X(x^* - \gamma \triangledown f(x^*))$$

Док-во:

$$\Pi_X(x^* - \gamma \triangledown f(x^*)) = \underset{x \in X}{argmin}\left[\|x - x^* + \gamma \triangledown f(x^*)\|_2^2\right]$$

$$= \underset{x \in X}{argmin}\left[\underbrace{\|x - x^*\|_2^2}_{\geq 0} + \underbrace{2\gamma \langle \triangledown f(x^*); x - x^* \rangle}_{? \geq 0 \text{ по условию оптимальности}} + \gamma^2 \|\triangledown f(x^*)\|_2^2\right]$$

$$\geq 0, \text{ но } 0 \text{ достигается на } x = x^* \quad \blacksquare$$

Док-во сходимости:

$$\|x^{k+1} - x^*\|_2^2 = \|\Pi_X\left(x^k - \gamma \triangledown f(x^k)\right) - x^*\|_2^2$$

Зе св-во $x^* = \ldots$

$$= \|\Pi_X\left(x^k - \gamma \triangledown f(x^k)\right) - \Pi_X\left(x^* - \gamma \triangledown f(x^*)\right)\|_2^2$$

2е св-во $\|\Pi(x_1) - \Pi(x_2)\|_2 \leq \|x_1 - x_2\|_2$

$$\leq \|x^k - x^* - \gamma \triangledown f(x^k) + \gamma \triangledown f(x^*)\|_2^2$$

$$= \|x^k - x^*\|_2^2 - 2\gamma \langle \triangledown f(x^k) - \triangledown f(x^*); x^k - x^* \rangle$$
$$+ \gamma^2 \|\triangledown f(x^k) - \triangledown f(x^*)\|_2^2$$

$\mu$-сильная выпуклость $\langle \rangle$ и $L$-гладкость для $\|\ \|_2^2$

$$\leq \|x^k - x^*\|_2^2 + 2\gamma \langle \triangledown f(x^*); x^k - x^* \rangle$$
$$- 2\gamma\left(\frac{\mu}{2}\|x^k - x^*\|_2^2 + f(x^k) - f(x^*)\right)$$

$$+ 2L\gamma^2 \left( f(x^k) - f(x^*) - \langle \nabla f(x^*); x^k - x^* \rangle \right)$$

$$= (1-\gamma\mu)\|x^k - x^*\|_2^2$$
$$+ 2\gamma(\gamma L - 1)\left( f(x^k) - f(x^*) - \langle \nabla f(x^*); x^k - x^* \rangle \right)$$

выпуклость ≥ 0
дивергенцией Брэгмана, кор. f

$$\gamma \leq \frac{1}{L}$$

$$\leq (1-\gamma\mu)\|x^k - x^*\|_2^2 \quad \blacksquare$$

Та же самая сходимость, что и у GD.

**Проекция:**

1) $X = \left\{ x \in \mathbb{R}^d \mid \|x\|_2^2 \leq 1 \right\}$ $\qquad \prod_X(x) = \min\left\{ 1; \frac{1}{\|x\|_2} \right\} x$



2) $X = \left\{ x \in \mathbb{R}^d \mid a_i \leq x_i \leq b_i \right\}$ $\qquad \left[ \prod_X(x) \right]_i = \begin{cases} a_i & x_i \leq a_i \\ x_i & a_i < x_i < b_i \\ b_i & x_i \geq b_i \end{cases}$

3) $X = \left\{ x \in \mathbb{R}^d \mid Ax = b \right\}$ $\qquad \prod_X(x) = x - A^\top (AA^\top)^{-1}(Ax - b)$

---

Линейная задача (как алгм. проекции/кв. задачи):

$$\min_{s \in X} \langle s; g \rangle$$

существует

1) $X = \{ x \in \mathbb{R}^d \mid \|x\|_1 \leq 1\}$

$$S^* = - \text{sign}(g_i) e_i \quad \leftarrow \text{базисный вектор}$$
$$i = \arg\max_j |g_j|$$

2) $X = \{ x \in \mathbb{R}^d \mid \sum_{i=1}^{d} x_i = 1 \; ; \; x_i \geq 0\}$

$$S^* = e_i \qquad i = \arg\min_j g_j$$

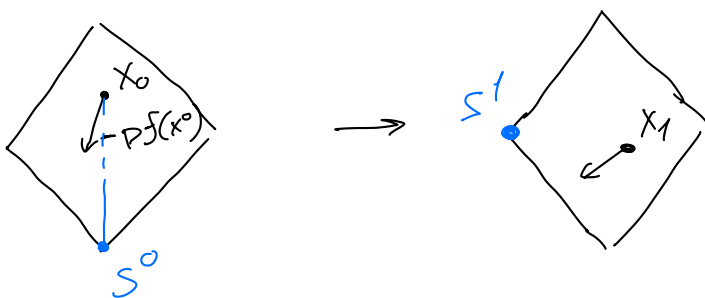3) $X = \{ x \in \mathbb{R}^d \mid \|x\|_\infty \leq 1\}$

$$S^* = - \sum_{i=1}^{d} \text{sign}(g_i) e_i$$

## Метод Франк - Вульфа

$$S^k = \arg\min_{s \in X} \langle s ; \nabla f(x^k) \rangle$$

$$x^{k+1} = (1-\gamma_k) x^k + \gamma_k S^k \qquad \gamma_k = \frac{2}{k+2}$$

## Физика:



$S^0, S^1 \ldots$ — на границе (в "углах" мн-ва)

$$x^{k+1} = \underbrace{\left(1 - \frac{2}{k+2}\right)}_{1-\gamma_k} x^k + \underbrace{\frac{2}{k+2}}_{\gamma_k} S^k$$

- смотрим на границу согласно лин.зад. $\Longleftarrow$
- усредняем точки на границах

$$x^{k+1} = \frac{k}{k+1} x^k + \frac{1}{k+1} S^k$$

подсчет среднего

<u>Док-во сходимости:</u>

$$f(x^{(k+1)}) = f(x^k + \gamma_k(s^k - x^k))$$

<span style="color:blue">$L$ — гладкость</span>

$$\leq f(x^k) + \gamma_k \langle s^k - x^k; \, \nabla f(x^k) \rangle + \frac{\gamma_k^2 L}{2} \| s^k - x^k \|_2^2$$

<span style="color:blue">$X$ — ограничено, $D = \text{diam} \, X$</span>

$$\leq f(x^k) + \gamma_k \langle s^k - x^k; \, \nabla f(x^k) \rangle + \frac{L D^2 \gamma_k^2}{2}$$

<span style="color:blue">$- f(x^*)$</span>

$$f(x^{(k+1)}) - f^* \leq f(x^k) - f^*$$
$$+ \gamma_k \langle s^k - x^k; \, \nabla f(x^k) \rangle + \frac{L D^2 \gamma_k^2}{2}$$

<span style="color:blue">$\langle s^k; \, \nabla f(x^k) \rangle = \min\limits_{s \in X} \langle s; \, \nabla f(x^k) \rangle \leq \langle x^*; \, \nabla f(x^k) \rangle$</span>

$$f(x^{(k+1)}) - f^* \leq f(x^k) - f^*$$
$$+ \gamma_k \langle x^* - x^k; \, \nabla f(x^k) \rangle + \frac{L D^2 \gamma_k^2}{2}$$

<span style="color:blue">выпуклость</span>

$$f(x^{(k+1)}) - f^* \leq (1 - \gamma_{\textcolor{red}{\boxed{k}}})(f(x^k) - f^*) + \frac{L D^2 \gamma_k^2}{2}$$

<span style="color:red">$\uparrow$ зависит $k$</span>

<span style="color:blue">По индукции: если $\gamma_k = \dfrac{2}{k+2}$, то

$$f(x^k) - f^* \leq \frac{\max\{4C; \, f(x^0) - f^*\}}{k+2}$$

$$C = \frac{L D^2}{2}$$</span>

$$f(x^{(k+1)}) - f^* \leq (1 - \gamma_k)(f(x^k) - f^*) + C \gamma_k^2$$

<span style="color:blue">$\overset{\gamma_k}{=} \left(1 - \frac{2}{k+2}\right)(f(x^k) - f^*) + \frac{4}{(k+2)^2} C$</span>

ПИ
$$\leq \left(1 - \frac{2}{k+2}\right)\left(\frac{\max\{4C;\ f(x^0)-f^*\}}{k+2}\right) + \frac{\max\{4C;\ f(x^0)-f^*\}}{(k+2)^2}$$

$$= \frac{\max\{4C;\ f(x^0)-f^*\}}{(k+1)+2}$$

Итоги ФВ:
- сублинейная схо. для выш. задачи (как у GD)
- в случае сильной выпуклости все равно сублинейно то рядок = линейная итер.