

Machine Learning: Optimization 4

Hiroshi Shimodaira and Hao Tang

2025

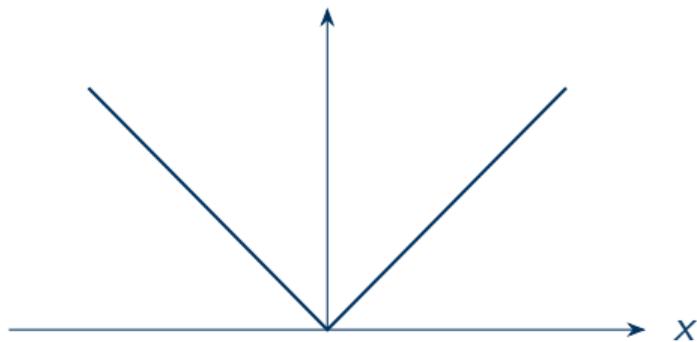
Ver. 1.0

Topics

- Subgradient
- Hinge loss
- Constrained optimisation problems
- Feasible solutions
- Lagrangian and Lagrange multiplier

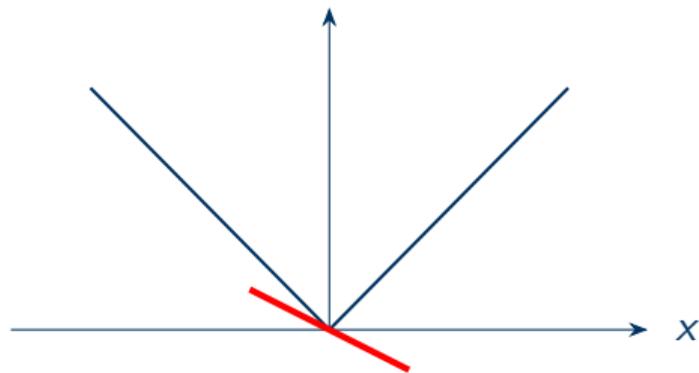
Subgradients for absolute values

$$f(x) = |x|$$



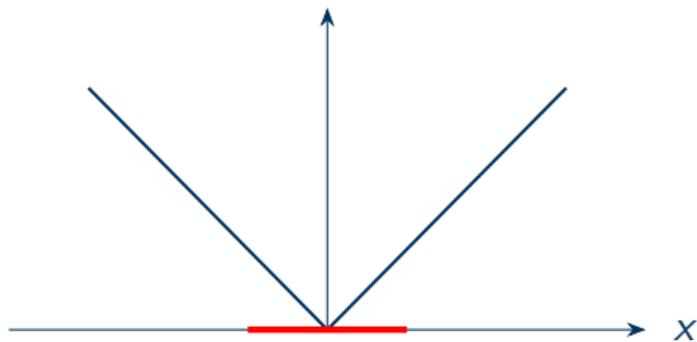
Subgradients for absolute values

$$f(x) = |x|$$



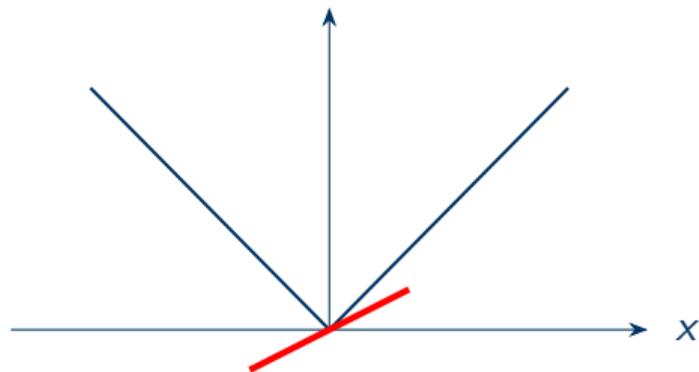
Subgradients for absolute values

$$f(x) = |x|$$



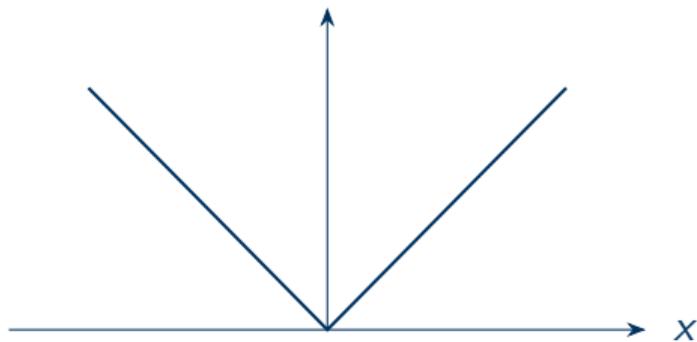
Subgradients for absolute values

$$f(x) = |x|$$

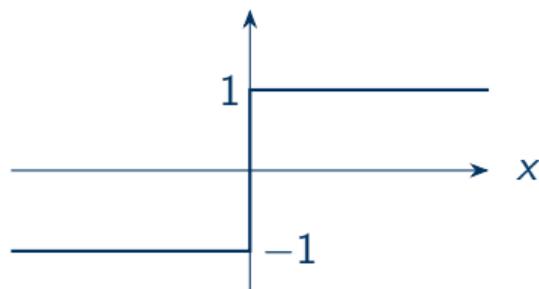


Subgradients for absolute values

$$f(x) = |x|$$



$$\partial|x|$$



$$\partial|x| = \begin{cases} \{-1\} & \text{if } x < 0 \\ [-1, 1] & \text{if } x = 0 \\ \{+1\} & \text{if } x > 0 \end{cases}$$

Subgradient

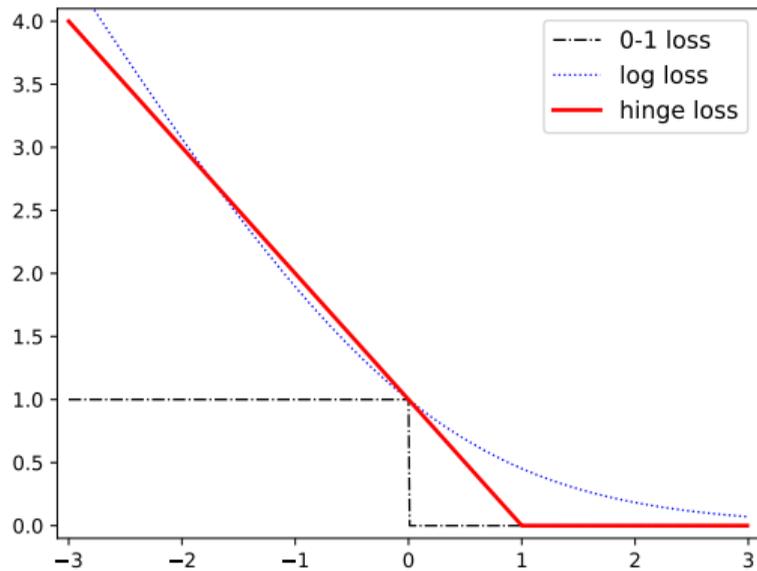
- A subgradient at \mathbf{x} is a vector \mathbf{g} that satisfies

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x}) \quad (1)$$

for any \mathbf{y} , and the set of subgradients at \mathbf{x} is denoted as $\partial f(\mathbf{x})$.

- Obviously, $\nabla f(\mathbf{x}) \in \partial f(\mathbf{x})$, if $\nabla f(\mathbf{x})$ exists.
- Convergence theorems can be ported to subgradient descent.

Hinge loss



Hinge loss (*cont.*)

- The hinge loss is defined as (\hat{y} : the raw output of classifier)

$$\ell_{\text{hinge}}(\hat{y}, y) = \max(0, 1 - \hat{y}y) \quad (2)$$

for a linear classifier

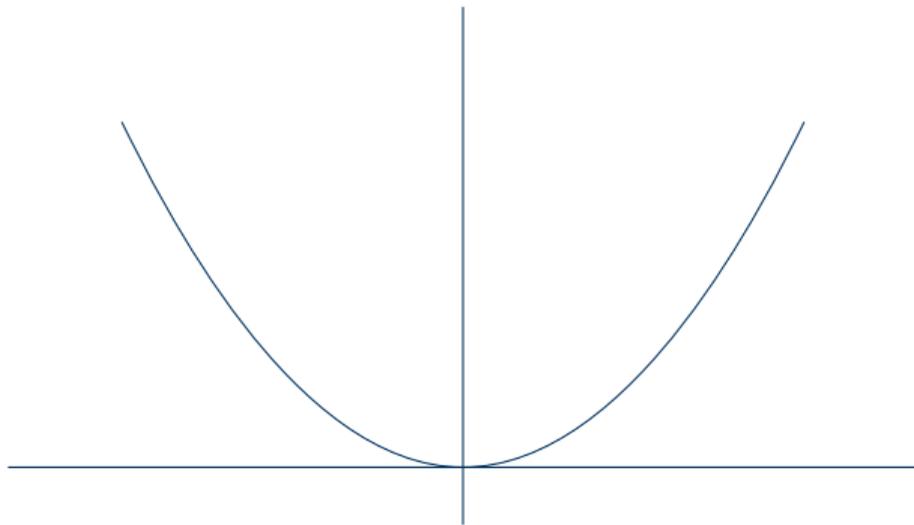
$$\ell_{\text{hinge}}(\mathbf{w}; \mathbf{x}, y) = \max(0, 1 - y\mathbf{w}^T \mathbf{x}). \quad (3)$$

- Just like the absolute value, the hinge loss is continuous and convex, but it is not differentiable.

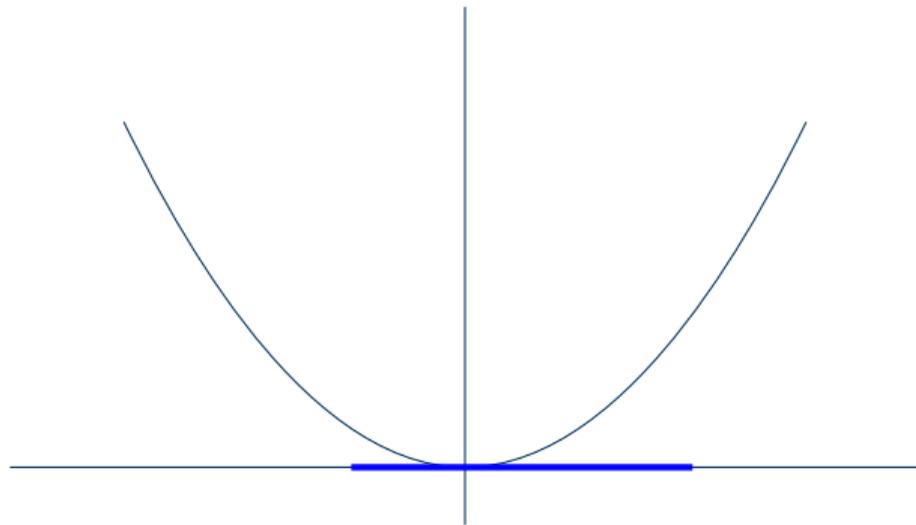
$$\nabla_{\mathbf{w}} \ell_{\text{hinge}} = \begin{cases} \mathbf{0} & \text{if } y\mathbf{w}^T \mathbf{x} \geq 1 \\ -y\mathbf{x} & \text{if } y\mathbf{w}^T \mathbf{x} < 1 \end{cases} \quad (4)$$

- When $y\mathbf{w}^T \mathbf{x} = 1$, we can pick and choose any vector that supports the loss function from below as the subgradient. In fact, $\mathbf{0}$ and $-y\mathbf{x}$ both work.

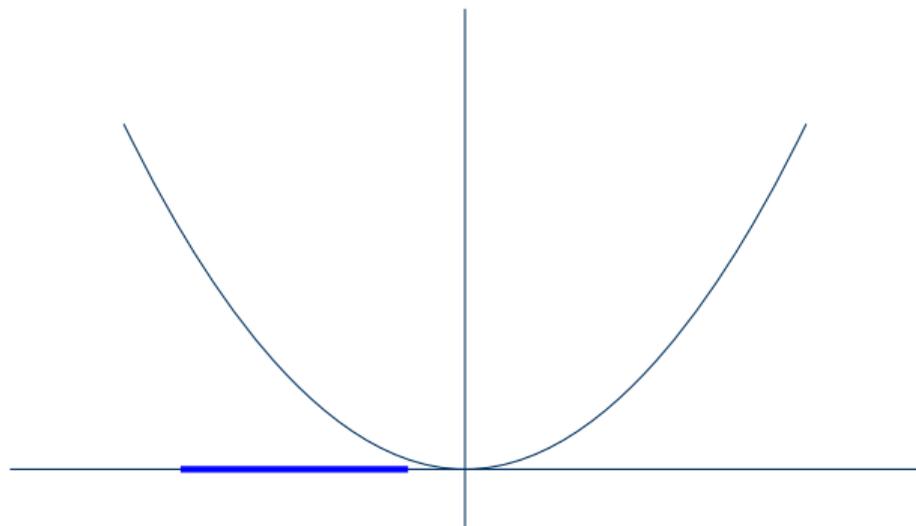
Constrained optimisation



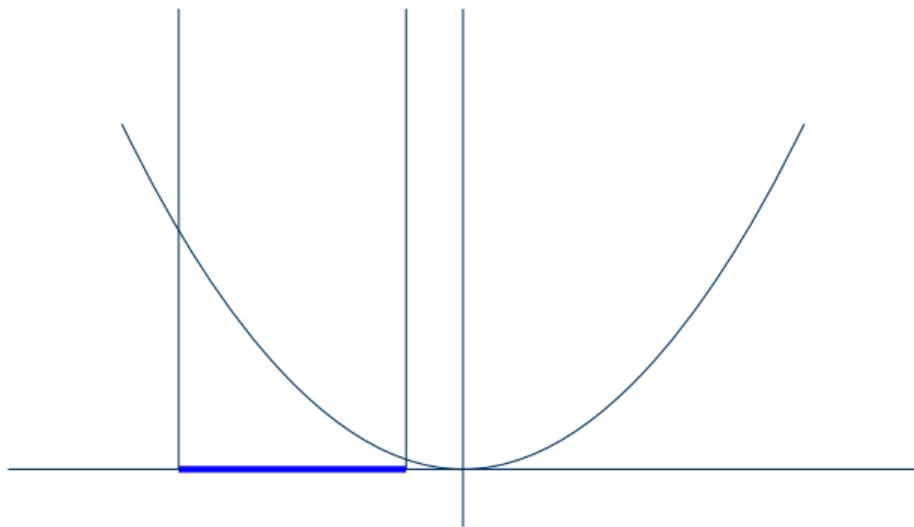
Constrained optimisation



Constrained optimisation



Setting up a barrier



An example optimisation-problem with constraints

- The problem

$$\begin{array}{ll} \min_x & x^2 \\ \text{s.t.} & -2.5 \leq x \leq -0.5 \end{array} \quad (5)$$

is an example of a constrained optimisation problem.

- The inequality $-2.5 \leq x \leq -0.5$ is called a *constraint*.
- Solutions that satisfy the constraints are called **feasible** solutions.

Setting up a barrier

- The problem

$$\begin{array}{ll} \min_x & x^2 \\ \text{s.t.} & -2.5 \leq x \leq -0.5 \end{array} \quad (6)$$

is equivalent to

$$\min_x x^2 + V_-(x) \quad (7)$$

where

$$V_-(x) = \begin{cases} 0 & \text{if } -2.5 \leq x \leq -0.5 \\ \infty & \text{otherwise} \end{cases} \quad (8)$$

An example optimisation-problem with constraints

- The problem

$$\begin{aligned} \min_{\mathbf{w}} \quad & L(\mathbf{w}) \\ \text{s.t.} \quad & \|\mathbf{w}\|_2^2 \leq 1 \end{aligned} \tag{9}$$

is an example of a constrained optimisation problem.

- The inequality $\|\mathbf{w}\|_2^2 \leq 1$ is called a *constraint*.
- Solutions that satisfy the constraints are called **feasible** solutions.

Setting up a barrier

- We can write the optimisation problem as

$$\min_{\mathbf{w}} L(\mathbf{w}) + V_-(\|\mathbf{w}\|_2^2 - 1), \quad (10)$$

where

$$V_-(s) = \begin{cases} 0 & \text{if } s \leq 0 \\ \infty & \text{if } s > 0 \end{cases}. \quad (11)$$

- This does not change anything; both problems are equally hard (or easy) to solve.

Soften the constraints

- We can approximate

$$\min_{\mathbf{w}} L(\mathbf{w}) + V_-(\|\mathbf{w}\|_2^2 - 1) \quad (12)$$

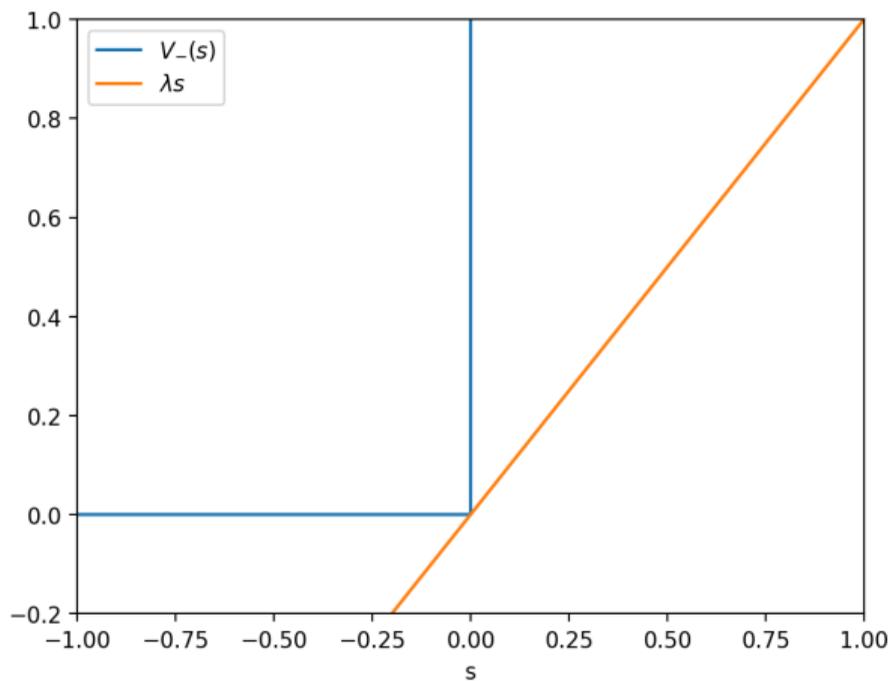
with

$$\min_{\mathbf{w}} L(\mathbf{w}) + \lambda(\|\mathbf{w}\|_2^2 - 1), \quad (13)$$

for some $\lambda \geq 0$.

- Note that $\lambda s \leq V_-(s)$ for all s .

Soften the constraints (*cont.*)



Lagrangian

- In general, if you have a optimisation problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & h(\mathbf{x}) \leq 0 \end{aligned} \tag{14}$$

the **Lagrangian** is defined as

$$f(\mathbf{x}) + \lambda h(\mathbf{x}) \tag{15}$$

for $\lambda \geq 0$.

- The value λ is called the **Lagrange multiplier**.

Solving the Lagrangian

- Solve $g(\lambda) = \min_{\mathbf{x}} [f(\mathbf{x}) + \lambda h(\mathbf{x})]$ for a particular λ .
- Find $\hat{\lambda}$ such that $\min_{\mathbf{x}} [f(\mathbf{x}) + \hat{\lambda} h(\mathbf{x})]$ gives a feasible solution.
- Suppose $\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} [f(\mathbf{x}) + \hat{\lambda} h(\mathbf{x})]$ and $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}: h(\mathbf{x}) \leq 0} f(\mathbf{x})$.

$$f(\hat{\mathbf{x}}) + \hat{\lambda} h(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*) + \hat{\lambda} h(\mathbf{x}^*) \leq f(\mathbf{x}^*) \quad (16)$$

Solving the Lagrangian (*cont.*)

- We want $f(\hat{\mathbf{x}}) = f(\hat{\mathbf{x}}) + \hat{\lambda}h(\hat{\mathbf{x}})$ leading to $f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*)$, so that we can conclude $f(\hat{\mathbf{x}}) = f(\mathbf{x}^*)$.
- If we want $\hat{\lambda}h(\hat{\mathbf{x}}) = 0$, then either $\hat{\lambda} = 0$ or $h(\hat{\mathbf{x}}) = 0$.
 - When $\hat{\lambda} = 0$, the minimiser of f is a feasible solution already.
 - When $h(\hat{\mathbf{x}}) = 0$, the minimiser of f is not necessarily a feasible solution, and we are on the edge of a constraint.

Example 1 - training of a word unigram model

Row, row, row your boat, gently down the stream
Merrily, merrily, merrily, merrily, life is but a dream

Example 1 - training of a word unigram model

Row, row, row your boat, gently down the stream
Merrily, merrily, merrily, merrily, life is but a dream

- There are 18 words.
- Intuitively,

$$p(\text{row}) = \frac{3}{18} \quad p(\text{merrily}) = \frac{4}{18} \quad p(\text{is}) = \frac{1}{18} \quad (17)$$

Example 1 - training of a word unigram model (*cont.*)

- There are 13 unique words.
- We refer to the set of unique words $V = \{\text{row, your, boat, gently, down, the, stream, merrily, life, is, but, a, dream}\}$ as the vocabulary.
- We assign each word v a probability β_v .
- The probability of a word is

$$p(w) = \prod_{v \in V} \beta_v^{\mathbb{1}_{v=w}}. \quad (18)$$

Example 1 - training of a word unigram model (*cont.*)

- We assume that each word is independent of others.
- This assumption is obviously wrong, but can go really far.
- The likelihood of β given the data is

$$\log p(w_1, \dots, w_N) = \log \prod_{i=1}^N p(w_i) = \log \prod_{i=1}^N \prod_{v \in V} \beta_v^{\mathbb{1}_{v=w_i}}. \quad (19)$$

- Since β is a probability vector, we have the assumption

$$\sum_{v \in V} \beta_v = 1. \quad (20)$$

Example 1 - training of a word unigram model (*cont.*)

- We arrive at the optimisation problem

$$\begin{aligned} \min_{\beta} \quad & - \sum_{i=1}^N \sum_{v \in V} \mathbb{1}_{v=w_i} \log \beta_v \\ \text{s.t.} \quad & \sum_{v \in V} \beta_v = 1 \end{aligned} \tag{21}$$

- Its Lagrangian is

$$F = - \sum_{i=1}^N \sum_{v \in V} \mathbb{1}_{v=w_i} \log \beta_v + \lambda \left(\sum_{v \in V} \beta_v - 1 \right). \tag{22}$$

Example 1 - training of a word unigram model (*cont.*)

- Solving the optimality condition gives

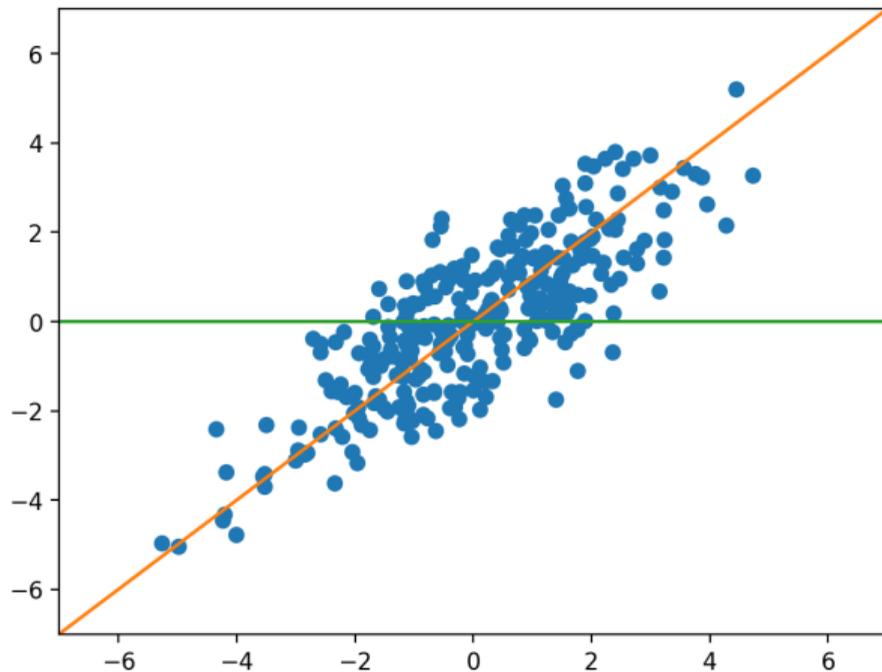
$$\frac{\partial F}{\partial \beta_k} = \sum_{i=1}^N \mathbb{1}_{k=w_i} \frac{1}{\beta_k} - \lambda = 0 \implies \beta_k = \frac{1}{\lambda} \sum_{i=1}^N \mathbb{1}_{k=w_i}. \quad (23)$$

Example 1 - training of a word unigram model (*cont.*)

$$\sum_{v \in \mathcal{V}} \beta_v = \sum_{v \in \mathcal{V}} \frac{1}{\lambda} \sum_{i=1}^N \mathbb{1}_{v=w_i} = 1 \implies \lambda = \sum_{v \in \mathcal{V}} \sum_{i=1}^N \mathbb{1}_{v=w_i} = N \quad (24)$$

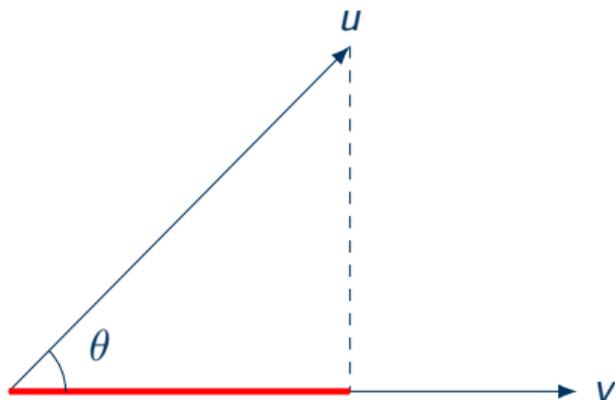
$$\beta_k = \frac{\sum_{i=1}^N \mathbb{1}_{k=w_i}}{\sum_{v \in \mathcal{V}} \sum_{i=1}^N \mathbb{1}_{v=w_i}} = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{k=w_i} \quad (25)$$

Example 2 - finding the best projection line/hyperplane



Projection of a vector

Projection of \mathbf{u} onto/from \mathbf{v}



$$\|\mathbf{u}\|_2 \cos \theta = \|\mathbf{u}\|_2 \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\|_2 \|\mathbf{v}\|_2} = \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{v}\|_2} \quad (26)$$

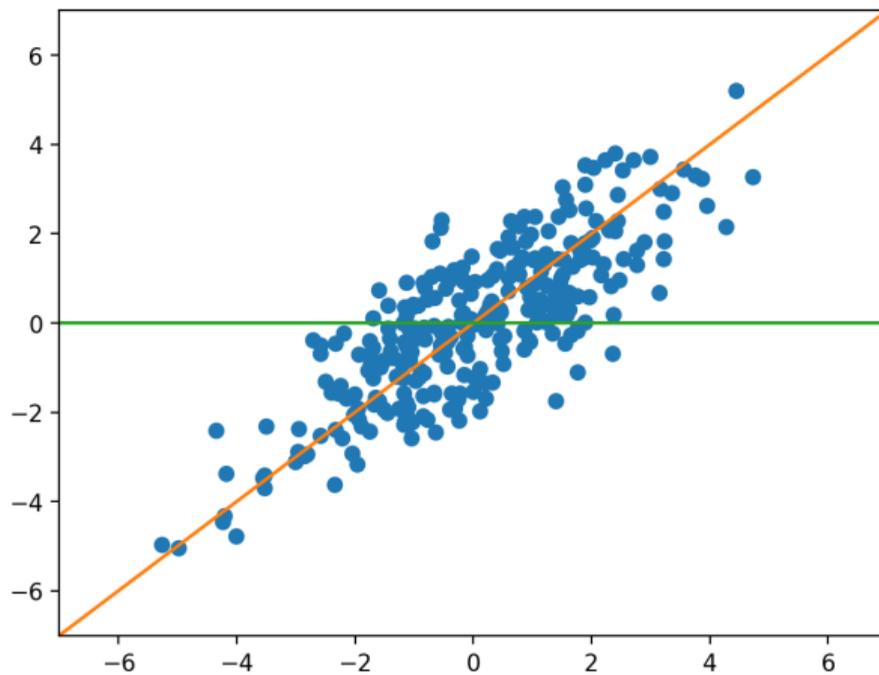
Example 2 - finding the best projection line/hyperplane (cont.)

- The projection of \mathbf{x} onto \mathbf{w} is $\frac{\mathbf{x}^\top \mathbf{w}}{\|\mathbf{w}\|_2}$.
- If we have N data points $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, then the sum of the (squared) projection is

$$\sum_{i=1}^N \left(\frac{|\mathbf{x}_i^\top \mathbf{w}|}{\|\mathbf{w}\|_2} \right)^2 = \frac{\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}{\mathbf{w}^\top \mathbf{w}}. \quad (27)$$

- The sum of squared projection can be seen as the spread of the data.

Maximal projection



Maximal projection (cont.)

- We want to find the maximum direction to project.
- The optimisation problem is

$$\max_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}}{\mathbf{w}^T \mathbf{w}}. \quad (28)$$

Maximal projection (cont.)

- The problem is scale invariant.

$$\frac{(a\mathbf{w})^\top \mathbf{X}^\top \mathbf{X} (a\mathbf{w})}{(a\mathbf{w})^\top (a\mathbf{w})} = \frac{\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}{\mathbf{w}^\top \mathbf{w}}. \quad (29)$$

- The problem is equivalent to

$$\max_{\mathbf{w}} \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} \quad \text{s.t.} \quad \|\mathbf{w}\|_2^2 = 1. \quad (30)$$

Maximal projection (cont.)

- The Lagrangian is

$$F = \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} + \lambda(1 - \|\mathbf{w}\|_2^2). \quad (31)$$

- Finding the optimal solution gives

$$\frac{\partial F}{\partial \mathbf{w}} = (\mathbf{X}^\top \mathbf{X} + \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\lambda \mathbf{w} = 0 \implies \mathbf{X}^\top \mathbf{X} \mathbf{w} = \lambda \mathbf{w}. \quad (32)$$

- It turns out that λ is an eigenvalue, and \mathbf{w} an eigenvector of $\mathbf{X}^\top \mathbf{X}$.

Maximal projection (cont.)

- Plugging the solution back to the objective,

$$\frac{\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}{\mathbf{w}^\top \mathbf{w}} = \frac{\lambda \mathbf{w}^\top \mathbf{w}}{\mathbf{w}^\top \mathbf{w}} = \lambda \quad (33)$$

- Since the goal is to find the maximal projection, this is now equivalent to finding the largest eigenvalue of $\mathbf{X}^\top \mathbf{X}$.

Maximal projection (cont.)

- The term

$$\frac{\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}{\mathbf{w}^\top \mathbf{w}} \quad (34)$$

is called the Rayleigh quotient.

- The optimal \mathbf{w} is called the first principal component.
- We will learn more about this when we talk about principal component analysis.

Quizzes

- Consider a set of two-dimensional data $\{\mathbf{x}_i\}_{i=1}^N$, where $\mathbf{x}_i = (x_{i1}, x_{i2})^\top$. Explain the difference between the best projection line (defined in the slides) and linear regression line from x_1 to x_2 (or from from x_2 to x_1).