

Minimum Mean Square Error Estimation

Logistic Regression

Minimum Mean Square Error Estimation

Minimum Mean Squared Error

recap . . .

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$$

$$\hat{y} = \arg \max_y p(y|\mathbf{x}) \quad \text{MAP estimate}$$

$$\hat{y} = \arg \max_y p(\mathbf{x}|y) \quad \text{ML (maximum likelihood) estimate}$$

- MAP minimizes risk under 0/1 loss.
- ML minimizes risk under 0/1 loss and uniform prior.

$$E[\ell(\hat{y}, y)] = E[\mathbb{1}_{\{\hat{y} \neq y\}}] = \mathbb{P}(\hat{y} \neq y)$$

MAP by definition picks the most probable y .

- MSE (mean squared error):

$$E[\ell(\hat{y}, y)] = E[(\hat{y} - y)^2]$$

- **claim:** minimum mean squared error estimate is given by

$$\hat{y} = f(\mathbf{x}) = E[y|\mathbf{x}] \quad \text{MMSE (minimum mean squared error estimate)}$$

(we need knowledge of $p(y|\mathbf{x})$ to find MMSE)

- note: MMSE is basis for denoising techniques such as Wiener filter, Kalman Filter.

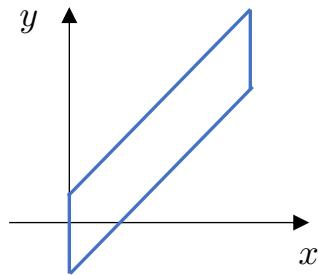
MMSE Example

- claim: the MMSE estimate of y given the features \mathbf{x} and complete knowledge of $p(\mathbf{x}, y)$ is:

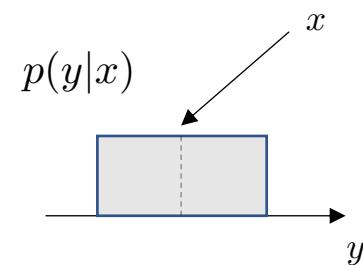
$$\hat{y} = f(\mathbf{x}) = E[y|\mathbf{x}]$$

- example:
 - given x , what is your best guess for y ?

$$p(x, y)$$



$$y|x \sim U[x - \Delta, x + \Delta]$$



- MMSE: $\hat{y} = E[y|x] = x$
- MAP: \hat{y} can be anything between $x - \Delta$ and $x + \Delta$

MMSE

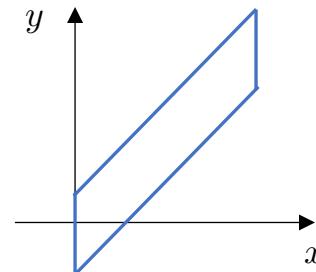
- claim: the MMSE estimate of y given the features \mathbf{x} and complete knowledge of $p(\mathbf{x}, y)$ is:

$$\hat{y} = f(\mathbf{x}) = E[y|\mathbf{x}] \quad \text{proof on canvas.}$$

- claim: the minimum mean squared error of the MMSE estimate of y given the features \mathbf{x} and $p(\mathbf{x}, y)$ is:

$$E[(y - E[y|\mathbf{x}])^2] = E_{\mathbf{x}}[\text{var}(y|\mathbf{x})]$$

$$p(x, y)$$



Logistic Regression

- logistic regression
- convex functions

Generative and Discriminative Classifiers

- if we have a good approximation of $p(\mathbf{x}, y)$, we know what to do:

$$\begin{aligned}\hat{y} &= \arg \max_y p(y|\mathbf{x}) && \text{MAP estimate} \\ \hat{y} &= \arg \max_y p(\mathbf{x}|y) && \text{ML (maximum likelihood) estimate} \\ \hat{y} &= E[y|\mathbf{x}] && \text{MMSE estimate}\end{aligned}$$

- probabilistic machine learning: we need to estimate $p(\mathbf{x}, y)$ from $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$
- most of the approaches we've studied have been *generative*
 - we imagine some distribution $p(\mathbf{x}|y)$ that *generates* our data
 - example: imagine that $p(\mathbf{x}|y) \sim \mathcal{N}(\mathbf{u}_i, \Sigma_i)$
- logistic regression is a *discriminative* classifier
 - discriminative classifiers directly find $p(y|\mathbf{x})$ to *discriminate* between classes.
 - directly fit $p(y|\mathbf{x})$ to data by treating \mathbf{x} as non-random
 - we don't learn anything about the distribution of x

(Binary) Logistic Regression

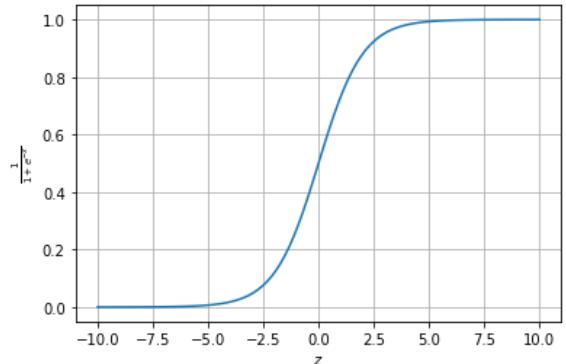
- logistic regression is a *discriminative* classifier
- we imagine that y is a coin flip with bias that depends on the features \mathbf{x}

$$y_i = \begin{cases} 1 & \text{with probability } \theta_i \\ -1 & \text{with probability } 1 - \theta_i \end{cases} \quad \theta_i = \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}}$$

$$p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{x}^T \mathbf{w}}} \quad p(y = -1|\mathbf{x}) = \frac{1}{1 + e^{\mathbf{x}^T \mathbf{w}}}$$

- testing: which is more probable?

$p(y = 1|\mathbf{x})$ or $p(y = -1|\mathbf{x})$ MAP estimate

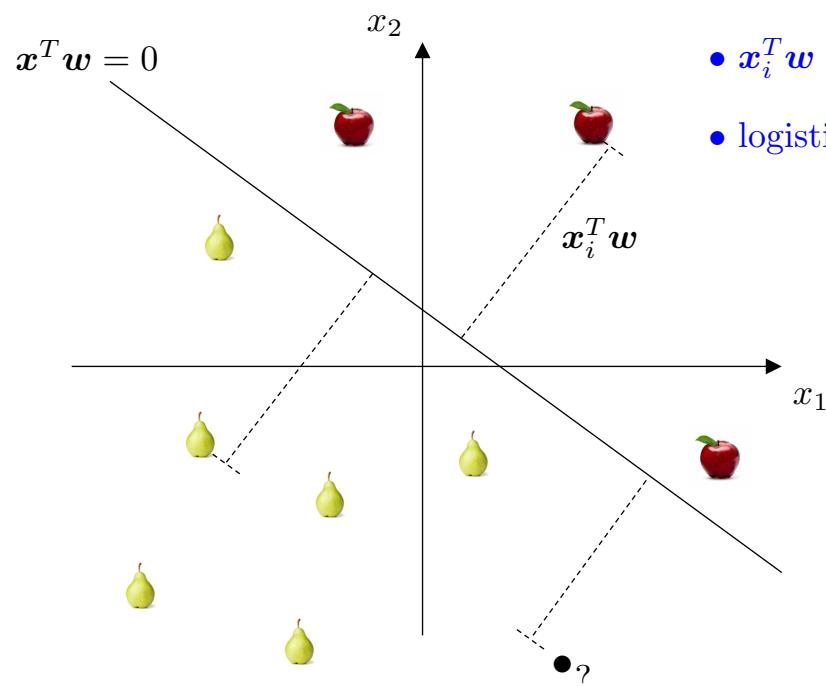


- decision boundary is $\mathbf{x}^T \mathbf{w} = 0$

$$\frac{1}{1 + e^{-z}} = \frac{1}{2} \quad \text{when } z = 0$$

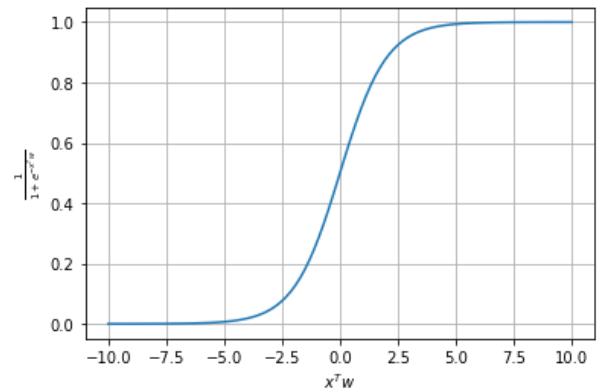
Logistic Regression

$$\mathbf{x}^T \mathbf{w} = [1 \ x_1 \ x_2] \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = w_0 + x_1 w_1 + x_2 w_2$$



- decision boundary is $\mathbf{x}^T \mathbf{w} = 0$
- $\mathbf{x}_i^T \mathbf{w}$ is the signed distance from the decision boundary
- logistic function takes $\mathbf{x}_i^T \mathbf{w}$ and squishes it to make it a probability

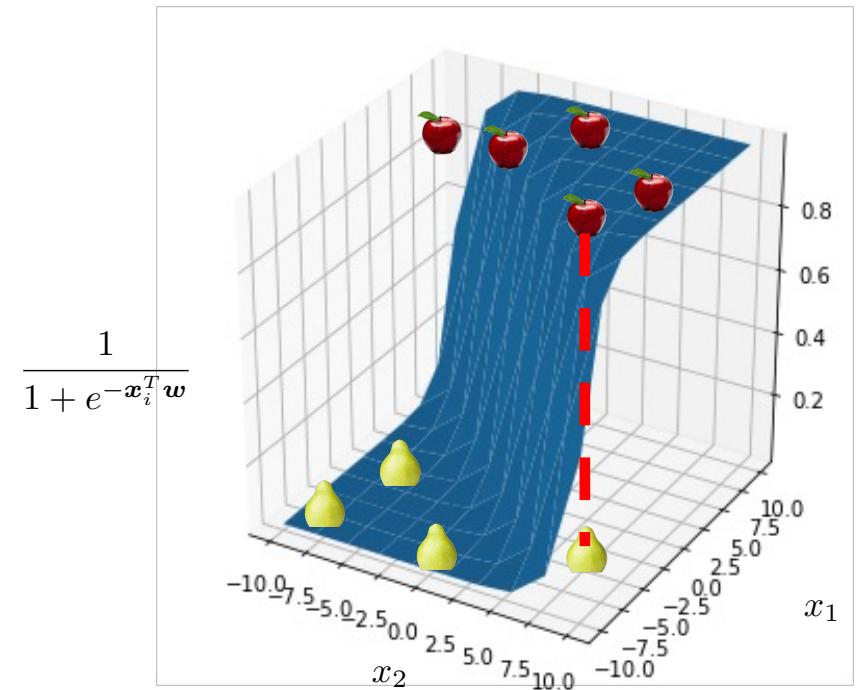
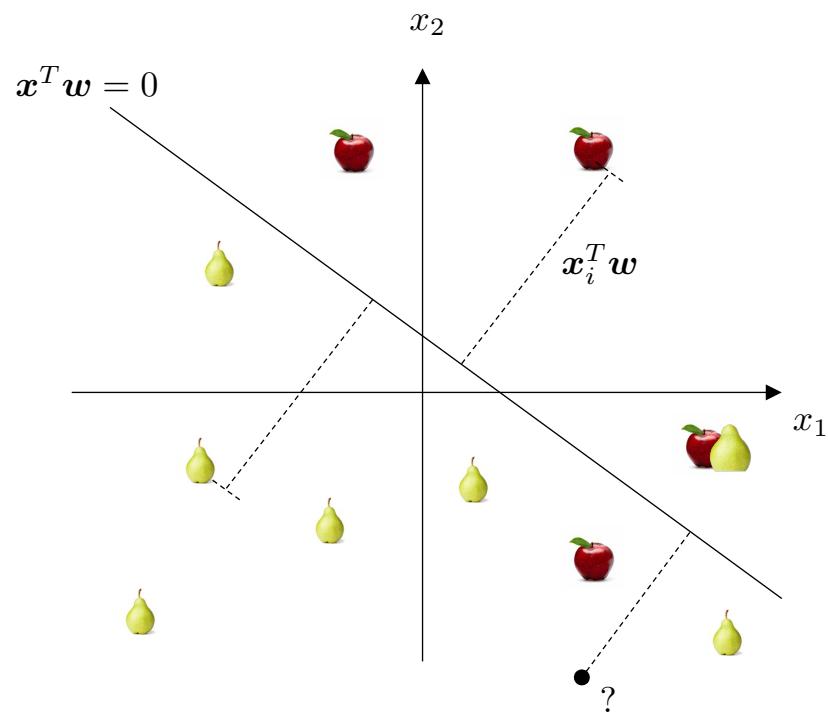
$$\frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}}$$



- probability that $?$ is a pear/apple depends distance from boundary

Logistic Regression

$$\mathbf{x}^T \mathbf{w} = [1 \ x_1 \ x_2] \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = w_0 + x_1 w_1 + x_2 w_2$$



Training



$y = -1$ if pear, $y = 1$ if apple

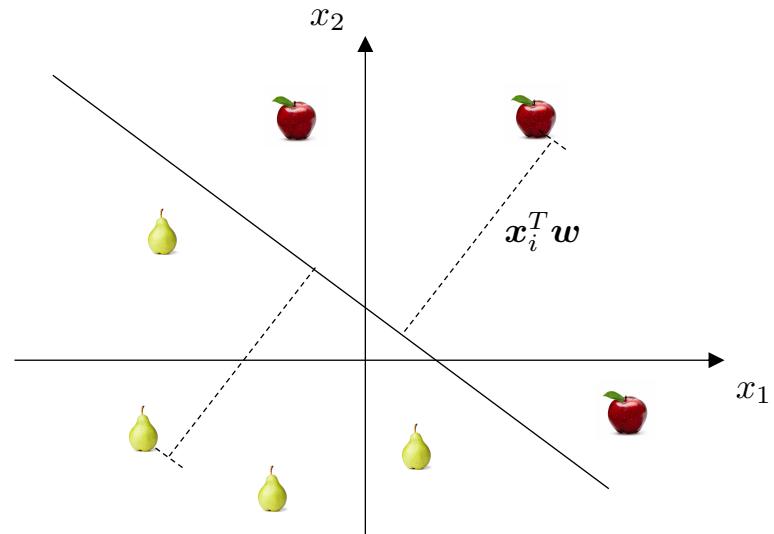
$$y_i = \begin{cases} 1 & \text{with probability } \theta_i \\ -1 & \text{with probability } 1 - \theta_i \end{cases}$$

$$\theta_i = \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}}$$

$$\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^m$$

- training: how do we find a good \mathbf{w} ?
 - maximum likelihood principle
 - given data (and our model) what's the most likely \mathbf{w} ?

$$p_{\mathbf{w}}(y_1, y_2, \dots, y_n) = \prod_i \theta_i^{\mathbb{I}\{y_i=1\}} (1 - \theta_i)^{\mathbb{I}\{y_i=-1\}}$$



Training

- maximum likelihood principle

$$p_{\mathbf{w}}(y_1, y_2, \dots, y_n) = \prod_i \theta_i^{\mathbb{I}\{y_i=1\}} (1 - \theta_i)^{\mathbb{I}\{y_i=-1\}}$$

$$\theta_i = \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}}$$

$$\log L(\mathbf{w}) = - \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{x}_i^T \mathbf{w}} \right)$$

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} \log L(\mathbf{w})$$

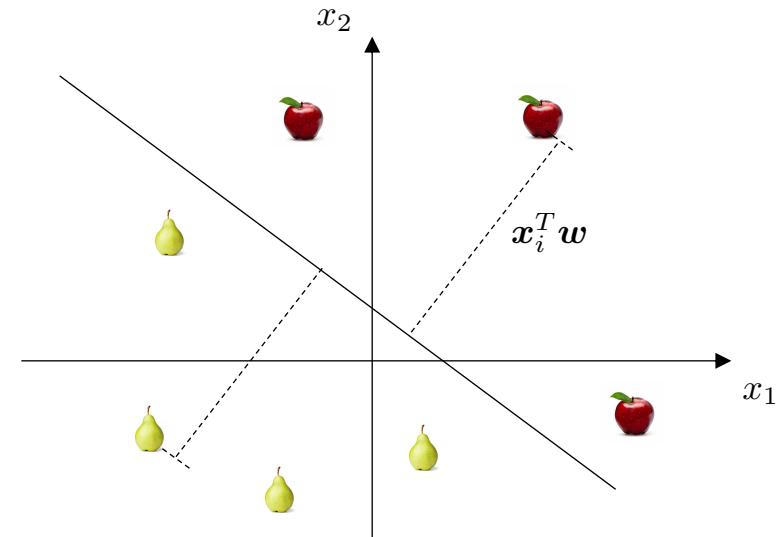
Training

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} - \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{x}_i^T \mathbf{w}} \right)$$

$$= \arg \min_{\mathbf{w}} \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{x}_i^T \mathbf{w}} \right)$$

- training a classifier in general:

$$\min_{\mathbf{w}} \sum_i \ell_i(\mathbf{w})$$



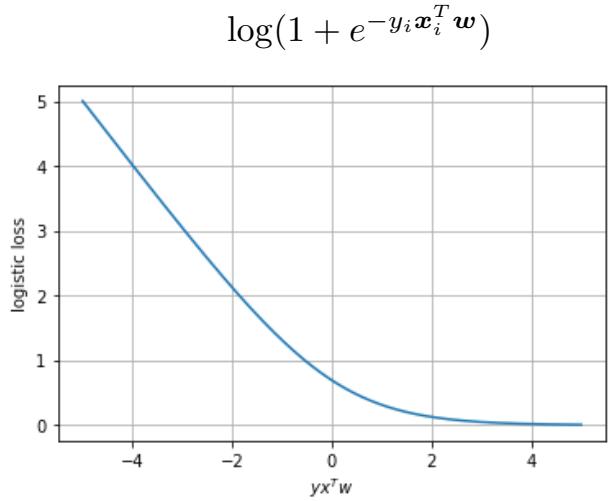
- hinge loss $(1 - y_i \mathbf{x}_i^T \mathbf{w})_+$
- squared error $(\mathbf{x}^T \mathbf{w} - y_i)^2$
- absolute loss $|y_i - \mathbf{x}_i^T \mathbf{w}|$
- logistic loss $\log(1 + e^{-y_i \mathbf{x}_i^T \mathbf{w}})$

Training

$$\begin{aligned}\mathbf{w}_{\text{ML}} &= \arg \max_{\mathbf{w}} - \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{x}_i^T \mathbf{w}} \right) \\ &= \arg \min_{\mathbf{w}} \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{x}_i^T \mathbf{w}} \right)\end{aligned}$$

- training a classifier in general:

$$\min_{\mathbf{w}} \sum_i \ell_i(\mathbf{w})$$



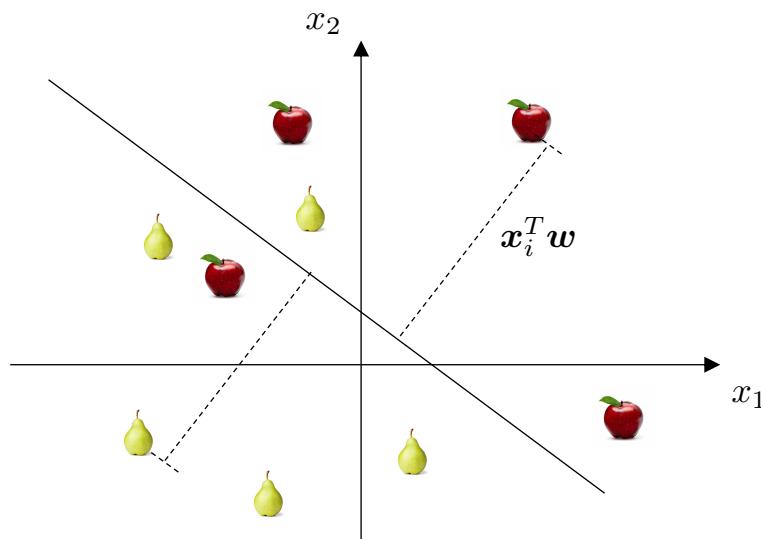
- hinge loss $(1 - y_i \mathbf{x}_i^T \mathbf{w})_+$
- squared error $(\mathbf{x}^T \mathbf{w} - y_i)^2$
- absolute loss $|y_i - \mathbf{x}_i^T \mathbf{w}|$
- logistic loss $\log(1 + e^{-y_i \mathbf{x}_i^T \mathbf{w}})$
- if $y_i \mathbf{x}^T \mathbf{w}$ is positive, loss is small
- if sign of $\mathbf{x}^T \mathbf{w}$ and y_i agree, $\log(1 + \epsilon) \approx 0$

Optimizing

- training reduces to solving the following optimization

$$\min_{\mathbf{w}} \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{x}_i^T \mathbf{w}} \right)$$

- changing \mathbf{w} change decision boundary orientation and steepness of logistic function



- no closed form solution
- iterative methods for finding \mathbf{w}
- logistic loss is easy to optimize . . .
- convex and differentiable

Convex Functions

- definition: a function $f(x)$ is convex if

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2) \quad \text{for all } \lambda \in [0, 1], \text{ all } x_1, x_2$$

- interpretation: the line connecting two points on the surface of the function is *always* above the function
- equivalently, a function is convex if the second derivative is non-negative

- examples:

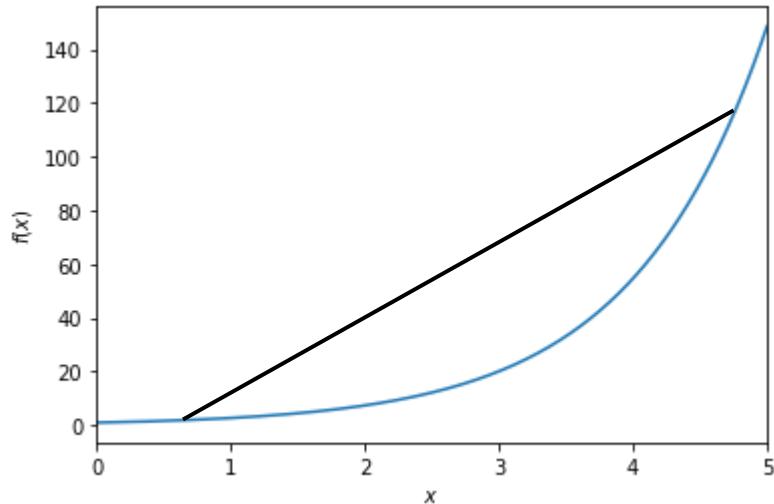
$$f(x) = x^2$$

$$f(x) = e^x$$

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

$$f(x) = \sin(x) \text{ for } \frac{\pi}{2} \leq x \leq \frac{3\pi}{2}$$

- a local minimum is a global minimum



Convex Functions in Multiple Dimensions

- definition: a function $f(\mathbf{x}), \mathbf{x} \in \mathbb{R}^n$ is convex if

$$f(\lambda\mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2) \leq \lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2) \text{ for all } \lambda \in [0, 1], \text{ all } \mathbf{x}_1, \mathbf{x}_2$$

- interpretation: line connecting any two points on the surface of the function is *always* above the function
- equivalently, a function is convex if the Hessian, $\nabla^2 f(\mathbf{x})$ is positive semi-definite

- examples: $f(\mathbf{x}) = \|\mathbf{x}\|_p \quad p \geq 1$

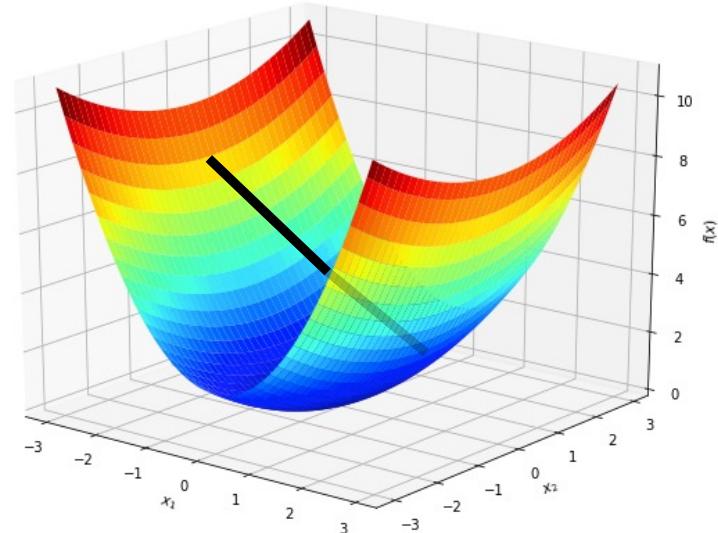
$$f(\mathbf{x}) = \max\{x_1, \dots, x_n\}$$

- important because computers can solve basically any optimization problem

$$\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$$

if $f(\mathbf{x})$ is a convex function over a convex set \mathcal{X}

- a local minimum is a global minimum



Optimizing

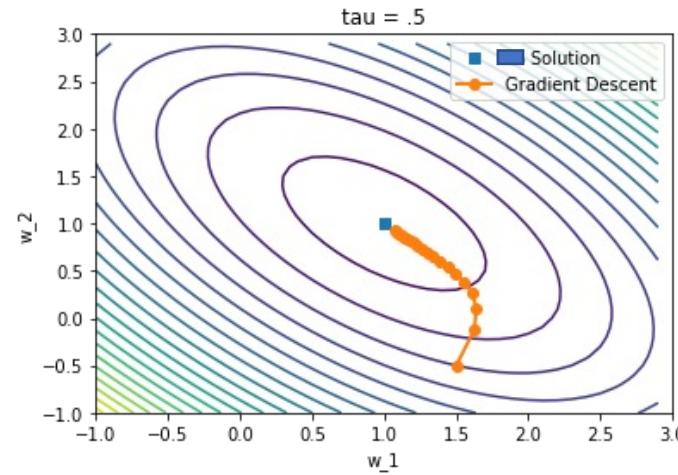
$$\min_{\mathbf{w}} \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{x}_i^T \mathbf{w}} \right)$$

- logistic loss is easy to optimize . . .
 - convex and differentiable
 - we can do gradient descent and are guaranteed to find the minimum

gradient descent

$$\text{for } k = 1 \dots \\ \mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \tau \nabla f(\mathbf{w})$$

$$\nabla_{\mathbf{w}} f(\mathbf{w}) = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix}$$



Finding the gradient

$$\min_{\mathbf{w}} \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{x}_i^T \mathbf{w}} \right) \quad \nabla_{\mathbf{w}} f(\mathbf{w}) = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix}$$

$$\nabla_{\mathbf{w}} f(\mathbf{w}) = \sum_i \frac{-y_i}{1 + e^{y_i \mathbf{w}^T \mathbf{x}_i}} \mathbf{x}_i$$

gradient descent

$$\text{for } k = 1 \dots \\ \mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \tau \nabla f(\mathbf{w})$$

Showing Convexity

$$\min_{\mathbf{w}} \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{x}_i^T \mathbf{w}} \right)$$

- option one: show Hessian is positive semi-definite
- option two: use some rules about checking convexity

$$1) \frac{d^2}{dz^2} \log(1 + e^{-z}) = \frac{e^z}{(e^z + 1)^2}$$

2) composition with a linear transform preserves convexity

[Boyd Vandenberghe, p.79]

3) non-negative sum of convex functions is convex

*See link BV on Canvas
to download the book*