

# CSCI 567: Machine Learning

Vatsal Sharan  
Fall 2022

Lecture 2, Sep 1

## Administrivia

- HW1 is out
- Due in about 2 weeks (9/14 at 2pm). **Start early!!!**
- Post on Ed Discussion if you're looking for teammates.

# Recap

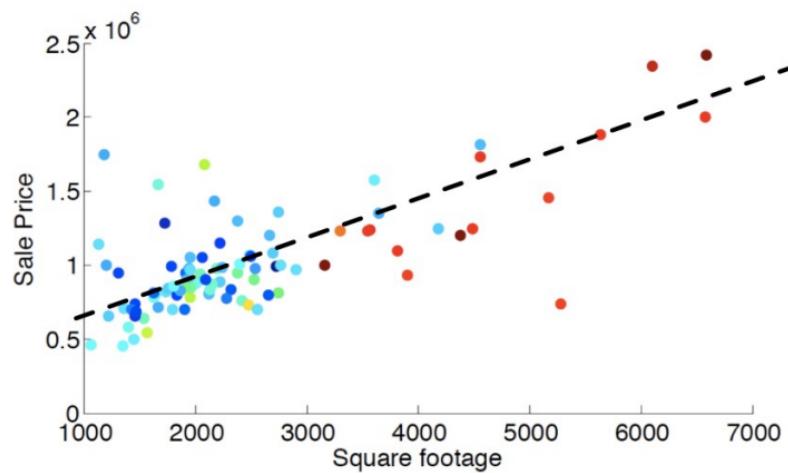
# Supervised learning in one slide

- Loss function:** What is the right loss function for the task?
- Representation:** What class of functions should we use?
- Optimization:** How can we efficiently solve the empirical risk minimization problem?
- Generalization:** Will the predictions of our model transfer gracefully to unseen examples?

*All related! And the fuel which powers everything is **data**.*

# Linear regression

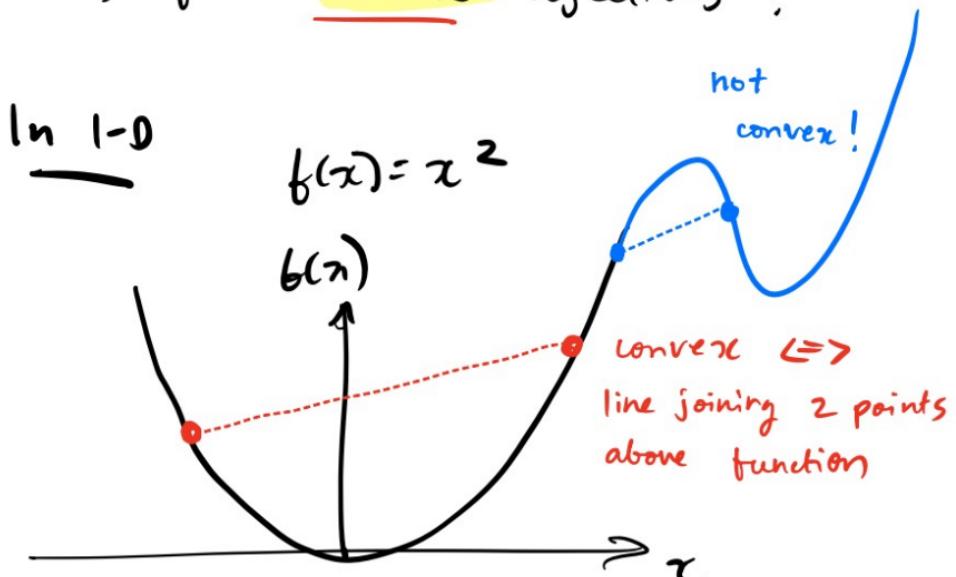
Predicted sale price = **price\_per\_sqft** × square footage + **fixed\_expense**



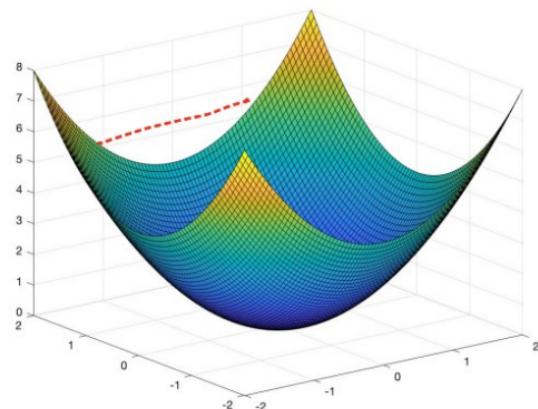
How to solve this? Find **stationary points**

## Are stationary points minimizers?

Yes, for convex objectives !



In high dimensions:



$\nabla^2(f(x))$  is positive  
semi-definite (psd)

# General least square solution

## Objective

$$\text{RSS}(\tilde{\boldsymbol{w}}) = \sum_i (\tilde{\boldsymbol{x}}_i^T \tilde{\boldsymbol{w}} - y_i)^2$$

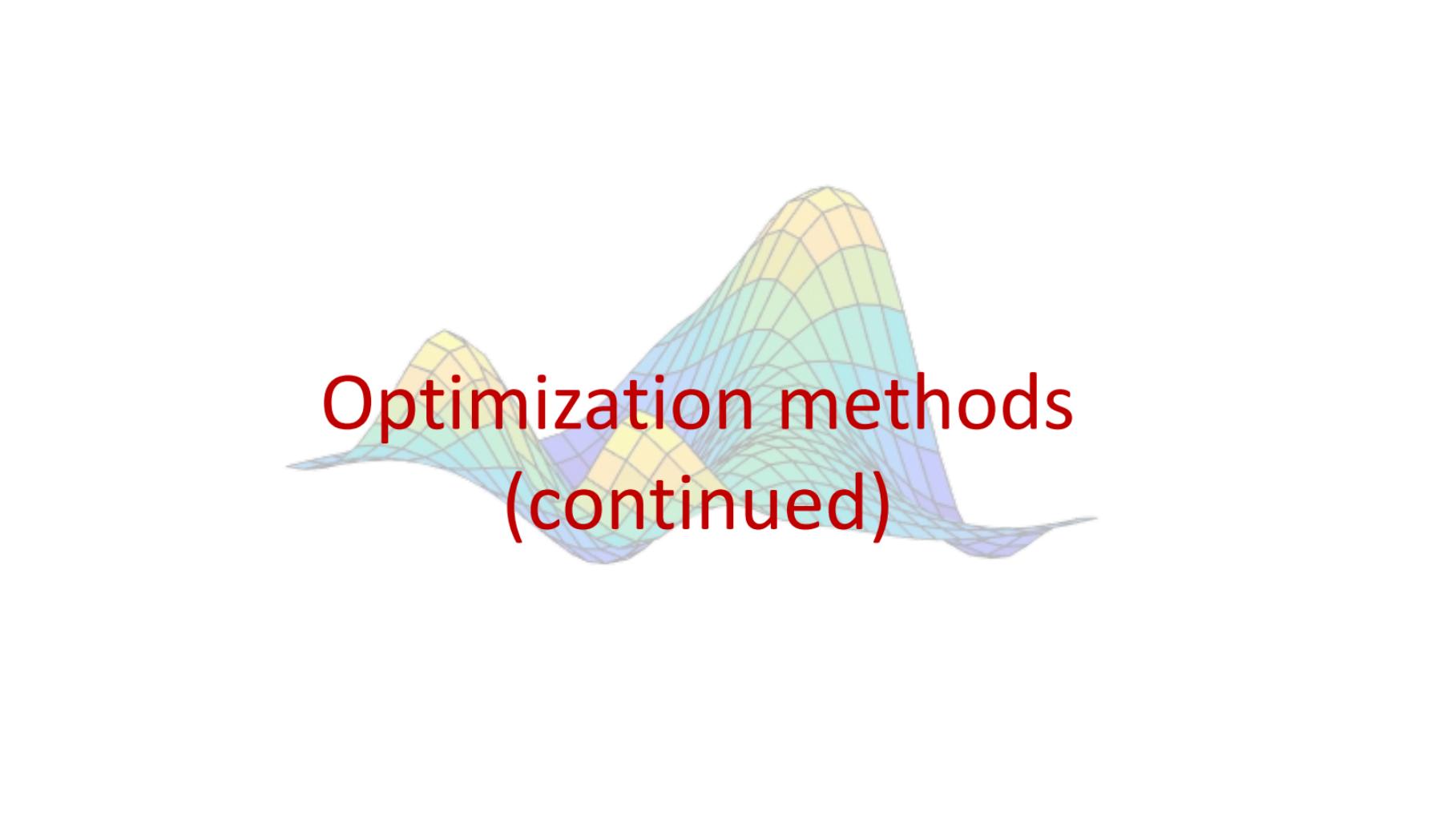
Find stationary points:

$$\begin{aligned}\nabla \text{RSS}(\tilde{\boldsymbol{w}}) &= 2 \sum_i \tilde{\boldsymbol{x}}_i (\tilde{\boldsymbol{x}}_i^T \tilde{\boldsymbol{w}} - y_i) \propto \left( \sum_i \tilde{\boldsymbol{x}}_i \tilde{\boldsymbol{x}}_i^T \right) \tilde{\boldsymbol{w}} - \sum_i \tilde{\boldsymbol{x}}_i y_i \\ &= (\tilde{\boldsymbol{X}}^T \tilde{\boldsymbol{X}}) \tilde{\boldsymbol{w}} - \tilde{\boldsymbol{X}}^T \boldsymbol{y}\end{aligned}$$

where

$$\tilde{\boldsymbol{X}} = \begin{pmatrix} \tilde{\boldsymbol{x}}_1^T \\ \tilde{\boldsymbol{x}}_2^T \\ \vdots \\ \tilde{\boldsymbol{x}}_n^T \end{pmatrix} \in \mathbb{R}^{n \times (d+1)}, \quad \boldsymbol{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$$

$$(\tilde{\boldsymbol{X}}^T \tilde{\boldsymbol{X}}) \tilde{\boldsymbol{w}} - \tilde{\boldsymbol{X}}^T \boldsymbol{y} = \mathbf{0} \quad \Rightarrow \quad \tilde{\boldsymbol{w}}^* = (\tilde{\boldsymbol{X}}^T \tilde{\boldsymbol{X}})^{-1} \tilde{\boldsymbol{X}}^T \boldsymbol{y}$$

A 3D surface plot of a bivariate normal distribution, showing a bell-shaped surface with a grid of lines. The surface transitions from yellow at the top to blue at the bottom. Overlaid on this surface is the text "Optimization methods (continued)" in a large, bold, red sans-serif font.

# Optimization methods (continued)

## Problem setup

Given: a function  $F(\mathbf{w})$

Goal: minimize  $F(\mathbf{w})$  (approximately)

Two simple yet extremely popular methods

**Gradient Descent (GD):** simple and fundamental

**Stochastic Gradient Descent (SGD):** faster, effective for large-scale problems

Gradient is the *first-order information* of a function.

Therefore, these methods are called *first-order methods*.

## Gradient descent

**GD**: keep moving in the *negative gradient direction*

Start from some  $\mathbf{w}^{(0)}$ . For  $t = 0, 1, 2, \dots$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)})$$

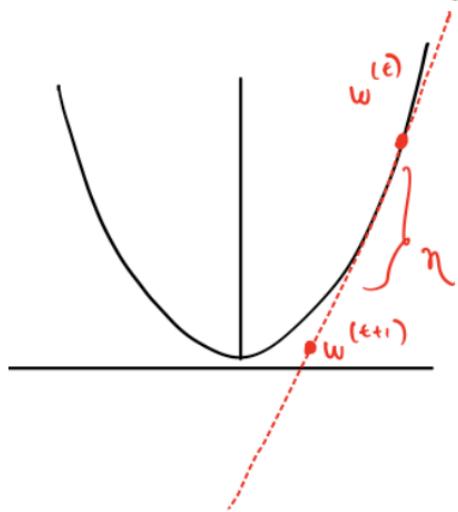
where  $\eta > 0$  is called step size or learning rate

- in theory  $\eta$  should be set in terms of some parameters of  $F$
- in practice we just try several small values
- might need to be changing over iterations (think  $F(w) = |w|$ )
- adaptive and automatic step size tuning is an active research area

# Why GD?

Intuition: First-order Taylor approximation

$$F(w) \approx F(w^{(t)}) + \nabla F(w^{(t)})^\top (w - w^{(t)})$$



$$\begin{aligned} F(w^{(t+1)}) &\approx F(w^{(t)}) - n \underbrace{\|\nabla F(w^{(t)})\|_2^2}_{>0} \\ \Rightarrow F(w^{(t+1)}) &\stackrel{\approx}{\leq} F(w^{(t)}) \end{aligned}$$

(this is only an approximation,  
and can be invalid if step  
size is too large)

## Convergence guarantees for GD

Many results for GD (and many variants) on *convex objectives*.

They tell you how many iterations  $t$  (in terms of  $\varepsilon$ ) are needed to achieve

$$F(\mathbf{w}^{(t)}) - F(\mathbf{w}^*) \leq \varepsilon$$

# Convergence guarantees for GD

Many results for GD (and many variants) on *convex objectives*.

They tell you how many iterations  $t$  (in terms of  $\varepsilon$ ) are needed to achieve

$$F(\mathbf{w}^{(t)}) - F(\mathbf{w}^*) \leq \varepsilon$$

Even for *nonconvex objectives*, some guarantees exist:

e.g. how many iterations  $t$  (in terms of  $\varepsilon$ ) are needed to achieve

$$\|\nabla F(\mathbf{w}^{(t)})\| \leq \varepsilon$$

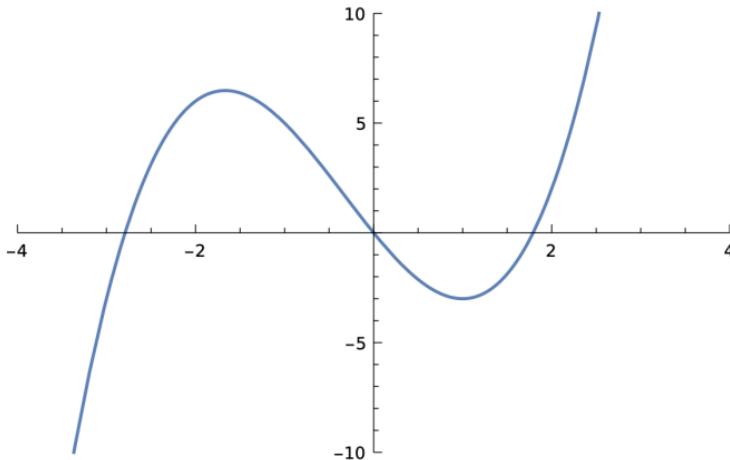
that is, how close is  $\mathbf{w}^{(t)}$  as an approximate stationary point

for convex objectives, stationary point  $\Rightarrow$  global minimizer

for nonconvex objectives, what does it mean?

## Stationary points: non-convex objectives

A stationary point can be a local minimizer or even a local/global maximizer (but the latter is not an issue for GD).

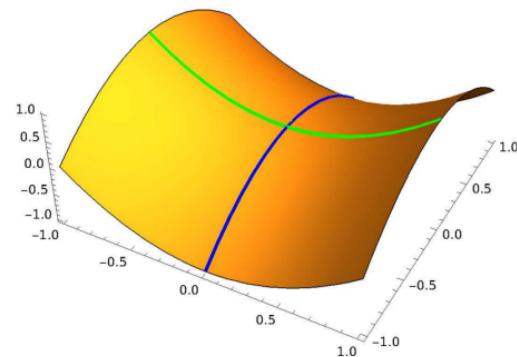


$$f(w) = w^3 + w^2 - 5w$$

# Stationary points: non convex objectives

A stationary point can also be *neither a local minimizer nor a local maximizer!*

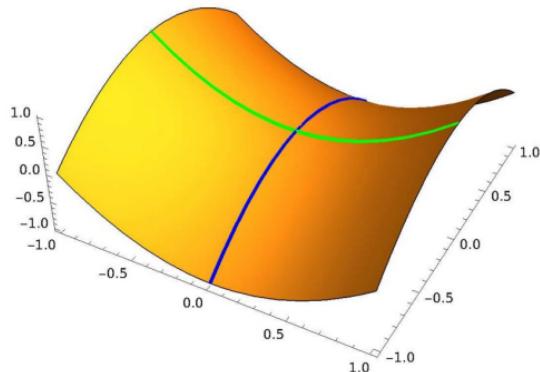
- $f(\mathbf{w}) = w_1^2 - w_2^2$
- $\nabla f(\mathbf{w}) = (2w_1, -2w_2)$
- so  $\mathbf{w} = (0, 0)$  is stationary
- local max for blue direction ( $w_1 = 0$ )
- local min for green direction ( $w_2 = 0$ )



Switch to Colab

# Stationary points: non convex objectives

This is known as a saddle point

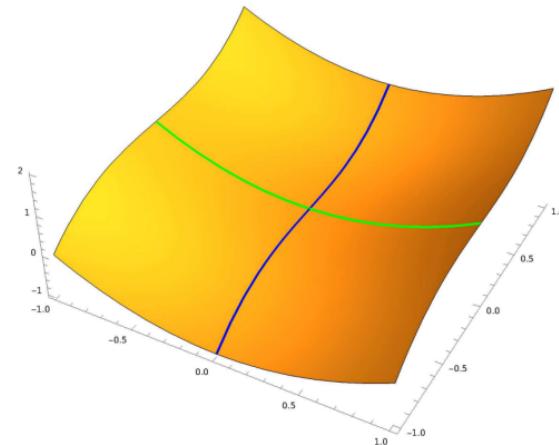


- but GD gets stuck at  $(0,0)$  only if initialized along the **green direction**
- so not a real issue especially *when initialized randomly*

# Stationary points: non convex objectives

But not all saddle points look like a “saddle” ...

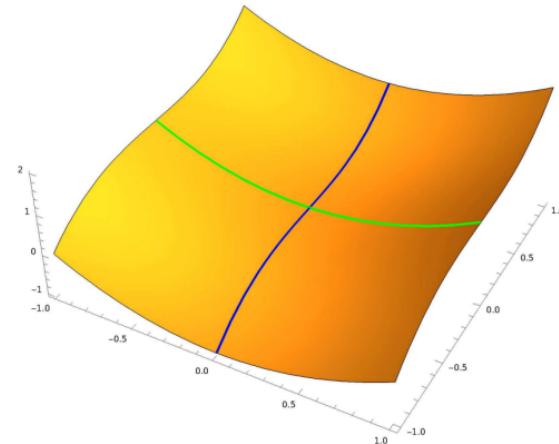
- $f(\mathbf{w}) = w_1^2 + w_2^3$
- $\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$
- so  $\mathbf{w} = (0, 0)$  is stationary
- not local min/max for blue direction  
( $w_1 = 0$ )



## Stationary points: non convex objectives

But not all saddle points look like a “saddle” ...

- $f(\mathbf{w}) = w_1^2 + w_2^3$
- $\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$
- so  $\mathbf{w} = (0, 0)$  is stationary
- not local min/max for blue direction ( $w_1 = 0$ )
- GD gets stuck at  $(0, 0)$  for *any initial point with  $w_2 \geq 0$  and small  $\eta$*



Even worse, distinguishing local min and saddle point is generally *NP-hard*.

# Stochastic Gradient descent

**GD**: keep moving in the *negative gradient direction*

**SGD**: keep moving in the *noisy negative gradient direction*

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \tilde{\nabla} F(\mathbf{w}^{(t)})$$

where  $\tilde{\nabla} F(\mathbf{w}^{(t)})$  is a random variable (called **stochastic gradient**) s.t.

$$\mathbb{E} [\tilde{\nabla} F(\mathbf{w}^{(t)})] = \nabla F(\mathbf{w}^{(t)}) \quad (\text{unbiasedness})$$

# Stochastic Gradient descent

**GD**: keep moving in the *negative gradient direction*

**SGD**: keep moving in the *noisy negative gradient direction*

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \tilde{\nabla} F(\mathbf{w}^{(t)})$$

where  $\tilde{\nabla} F(\mathbf{w}^{(t)})$  is a random variable (called **stochastic gradient**) s.t.

$$\mathbb{E} [\tilde{\nabla} F(\mathbf{w}^{(t)})] = \nabla F(\mathbf{w}^{(t)}) \quad (\text{unbiasedness})$$

- Key point: it could be much faster to obtain a stochastic gradient!
- Similar convergence guarantees, usually needs more iterations but each iteration takes less time.

## Summary: Gradient descent & Stochastic Gradient descent

- GD/SGD converges to a stationary point
- for convex objectives, this is all we need

## Summary: Gradient descent & Stochastic Gradient descent

- GD/SGD converges to a stationary point
- for convex objectives, this is all we need
- for nonconvex objectives, can get stuck at local minimizers or “bad” saddle points (random initialization escapes “good” saddle points)
- recent research shows that *many problems have no “bad” saddle points or even “bad” local minimizers*
- justify the practical effectiveness of GD/SGD (default method to try)

## Second-order methods

GD: first-order Taylor approximation

$$F(w) \approx F(w^{(t)}) + \nabla F(w^{(t)})^T (w - w^{(t)})$$

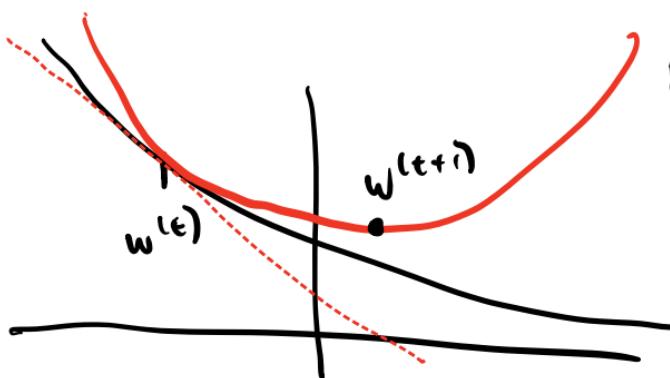
$$f(y) \approx f(x) + f'(x)(y-x) + \frac{f''(x)}{2}(y-x)^2$$

What about a second-order Taylor approximation?

$$F(w) \approx F(w^{(t)}) + \nabla F(w^{(t)})^T (w - w^{(t)}) + \frac{1}{2} (w - w^{(t)})^T H_t (w - w^{(t)})$$

where  $H_t = \nabla^2 F(w^{(t)}) \in \mathbb{R}^{d \times d}$  is Hessian of  $F$  at  $w^{(t)}$

$$H_{t,i,j} = \left. \frac{\partial^2 F(w)}{\partial w_i \partial w_j} \right|_{w=w^{(t)}}$$



Define  $\tilde{F}(w)$  := 2nd order approximation  
Set  $\nabla \tilde{F}(w) = 0$  to find minima.

$$\nabla F(w^{(t)}) + H_t w - \frac{H_t w^{(t)}}{2} - \frac{1}{2} H_t w^{(t)} = 0$$

$$H_t w = H_t w^{(t)} - \nabla F(w^{(t)})$$

$$\Rightarrow w = w^{(t)} - H_t^{-1} \nabla F(w^{(t)})$$

Newton's method:  $w^{(t+1)} \leftarrow w^{(t)} - H_t^{-1} \nabla F(w^{(t)})$

GD:  $w^{(t+1)} \leftarrow w^{(t)} - \eta \nabla F(w^{(t)})$

## Newton's method

- \* no learning rate
- \* super fast convergence
- \* Know & invert Hessian  
(inversion takes  $O(d^3)$   
time naively)

## GD

- \* need to tune  $\eta$
- \* slower convergence
- \* fast!  
 $O(d)$  time

# Linear classifiers



# The Setup

Recall the setup:

- input (feature vector):  $x \in \mathbb{R}^d$
- output (label):  $y \in [C] = \{1, 2, \dots, C\}$
- goal: learn a mapping  $f : \mathbb{R}^d \rightarrow [C]$

This lecture: **binary classification**

- Number of classes:  $C = 2$
- Labels:  $\{-1, +1\}$  (cat or dog)

## Representation: Choosing the **function class**

Let's follow the recipe, and pick a function class  $\mathcal{F}$ .

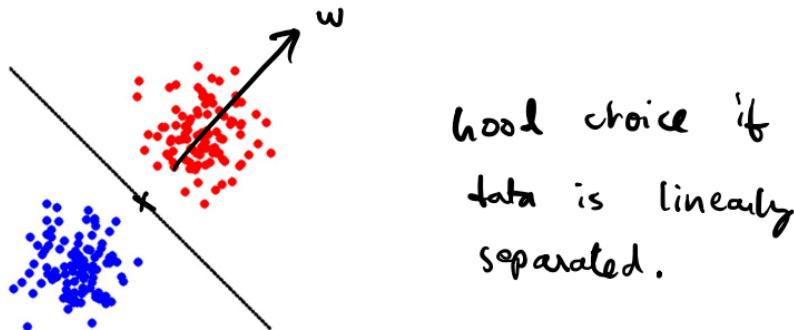
We continue with linear models, but how to predict a label using  $\mathbf{w}^T \mathbf{x}$ ?

*Sign* of  $\mathbf{w}^T \mathbf{x}$  predicts the label:

$$\text{sign}(\mathbf{w}^T \mathbf{x}) = \begin{cases} +1 & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ -1 & \text{if } \mathbf{w}^T \mathbf{x} \leq 0 \end{cases}$$

(Sometimes use sgn for sign too.)

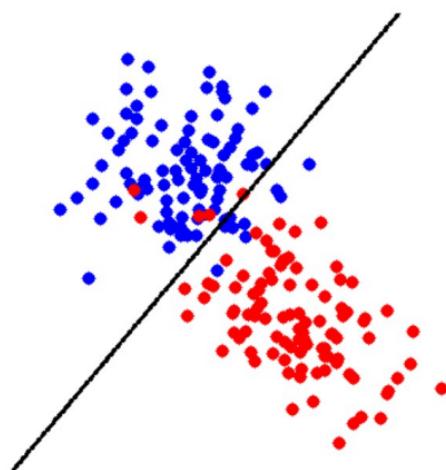
## Representation: Choosing the **function class**



Definition: The function class of separating hyperplanes is defined as:

$$\mathcal{F} = \{ f(x) = \text{sgn}(w^T x) : w \in \mathbb{R}^d \}.$$

Still makes sense for “almost” linearly separable data



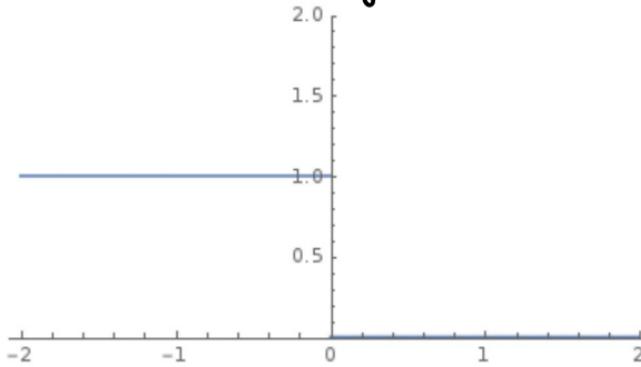
## Choosing the loss function

Most common loss  $l(f(x), y) = \mathbb{1}(f(x) \neq y)$

Loss as a function of  $yw^T x$

$$l_{0-1}(yw^T x) = \mathbb{1}(yw^T x \leq 0)$$

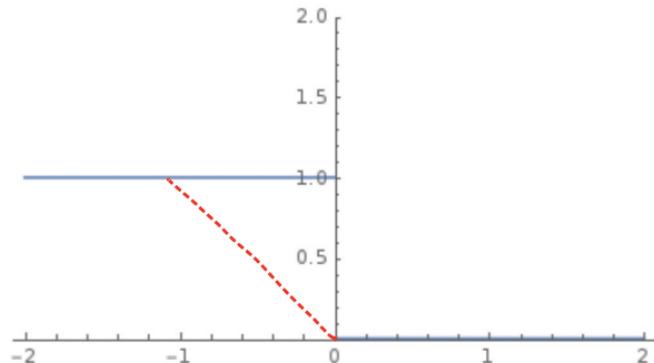
$$l_{0-1}(yw^T x)$$



$$yw^T x$$

## Choosing the loss function: **minimizing 0/1 loss is hard**

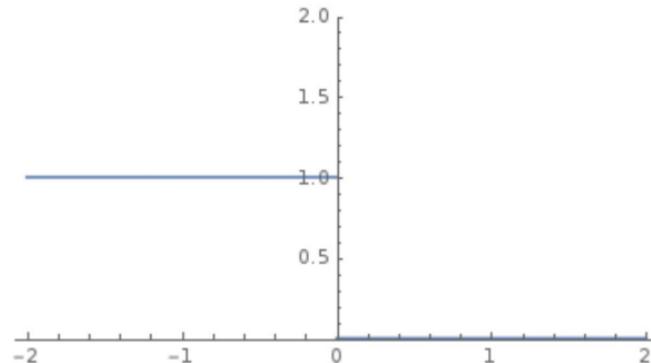
However, 0-1 loss is not convex.



Even worse, minimizing 0-1 loss is NP-hard in general.

## Choosing the loss function: **surrogate losses**

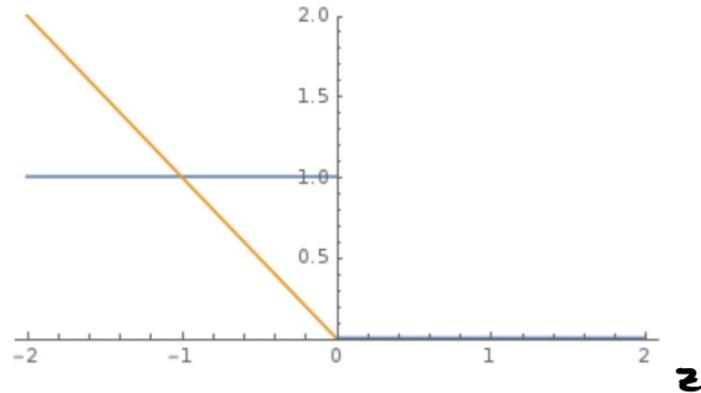
Solution: use a **convex surrogate loss**



## Choosing the loss function: **surrogate losses**

Solution: use a **convex surrogate loss**

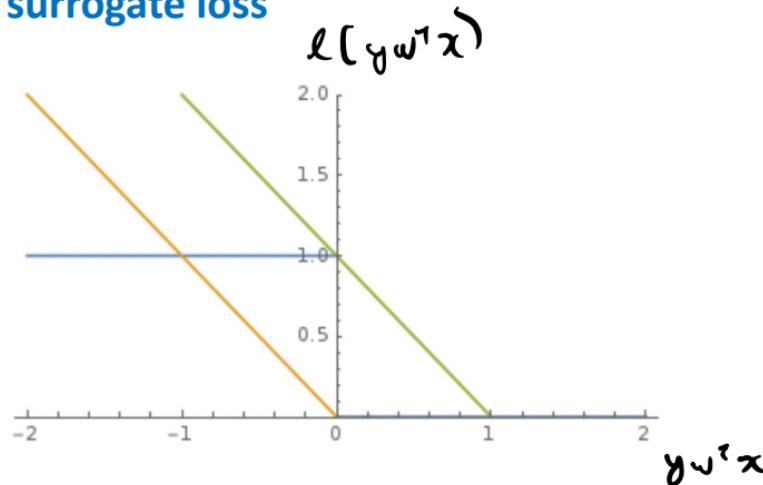
$$\ell(z)$$



- **perceptron loss**  $\ell_{\text{perceptron}}(z) = \max\{0, -z\}$  (used in Perceptron)

## Choosing the loss function: **surrogate losses**

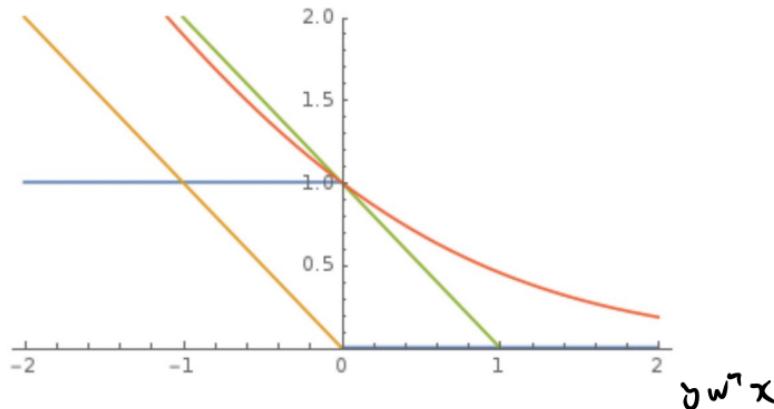
Solution: use a **convex surrogate loss**



- **perceptron loss**  $\ell_{\text{perceptron}}(z) = \max\{0, -z\}$  (used in Perceptron)
- **hinge loss**  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$  (used in SVM and many others)

# Choosing the loss function: **surrogate losses**

Solution: use a **convex surrogate loss**



- **perceptron loss**  $\ell_{\text{perceptron}}(z) = \max\{0, -z\}$  (used in Perceptron)
- **hinge loss**  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$  (used in SVM and many others)
- **logistic loss**  $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z))$  (used in logistic regression; the base of log doesn't matter)

## Onto Optimization!

Find ERM :

$$w^* = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \quad \frac{1}{n} \left( \sum_{i=1}^n l(y_i w^T x_i) \right)$$

where  $l(\cdot)$  is a convex surrogate loss.

- No closed-form solution in general (in contrast to linear regression)
- We can use our **optimization** toolbox!

*New York Times, 1958*

## NEW NAVY DEVICE LEARNS BY DOING

Psychologist Shows Embryo  
of Computer Designed to  
Read and Grow Wiser

# Perceptron

WASHINGTON, July 7 (UPI)—The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.

*The Navy last week demonstrated the embryo of an electronic computer named the **Perceptron** which, when completed in about a year, is expected to be the first non-living mechanism able to "perceive, recognize and identify its surroundings without human training or control."*

*New York Times, 1958*

## NEW NAVY DEVICE LEARNS BY DOING

Psychologist Shows Embryo  
of Computer Designed to  
Read and Grow Wiser

WASHINGTON, July 7 (UPI)—The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.

*The Navy last week demonstrated the embryo of an electronic computer named the **Perceptron** which, when completed in about a year, is expected to be the first non-living mechanism able to "perceive, recognize and identify its surroundings without human training or control."*

## Recall perceptron loss

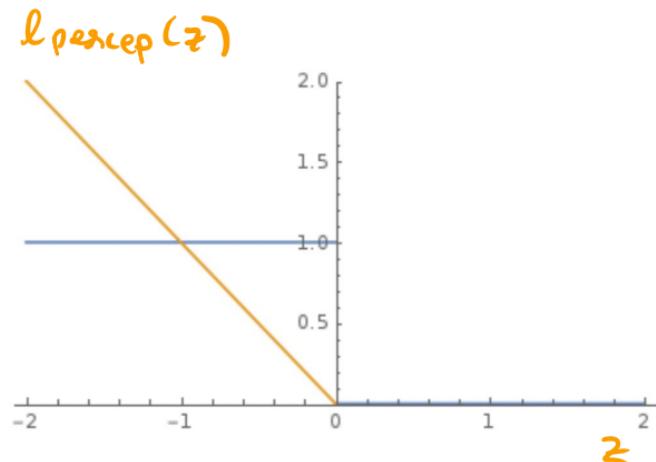
$$\begin{aligned} F(w) &= \frac{1}{n} \sum_{i=1}^n l_{\text{percep}}(y_i w^\top x_i) \\ &= \frac{1}{n} \sum_{i=1}^n \max \{0, -y_i w^\top x_i\} \end{aligned}$$

Let's try GD / SGD.

What's the gradient:

$$0, z > 0$$

$$-1, z \leq 0$$



## Applying GD to perceptron loss

Gradient is

$$\nabla F(w) = (1/n) \sum_{i=1}^n -\mathbb{1}[y_i w^\top x_i \leq 0] y_i x_i$$

(only misclassified examples count)

$$GD : w \leftarrow w + \underbrace{\frac{n}{n} \sum_{i=1}^n \mathbb{1}[y_i w^\top x_i \leq 0] y_i x_i}_{\text{need the entire training set for every GD update.}}$$

need the entire training set for every GD update.

## Applying SGD to perceptron loss

How to get a stochastic gradient?

→ pick one example  $i \in [n]$  uniformly at random, let apply

$$\tilde{\nabla} F(w^{(t)}) = -\mathbb{1}[y_i w^\top x_i \leq 0] y_i x_i$$

Unbiased. Why?  $\mathbb{E}[\tilde{\nabla} F(w^{(t)})] = \frac{1}{n} \sum_{i=1}^n -\mathbb{1}[y_i w^\top x_i \leq 0] y_i x_i = \nabla F(w^{(t)})$

SGD update:  $w \leftarrow w + \eta \mathbb{1}(y_i w^\top x_i \leq 0) y_i x_i$   
fast! one datapoint per update

Objective function of most ML tasks is a finite sum, trick generally

## Perceptron algorithm

Start with  $\omega = 0$  on perceptron loss.

Initialize  $\omega = 0$

Repeat

- Pick  $x_i \sim \text{Unif}(x_1, \dots, x_n)$
- If  $\text{sgn}(\omega^T x_i) \neq y_i$

$$\omega \leftarrow \omega + y_i x_i$$

## Perceptron algorithm: Intuition

Say that  $w$  makes mistake on  $(x_i, y_i)$

$$y_i w^T x_i < 0$$

Consider  $w' = w + y_i x_i$

$$y_i (w')^T x_i = y_i w^T x_i + y_i^2 x_i^T x_i$$

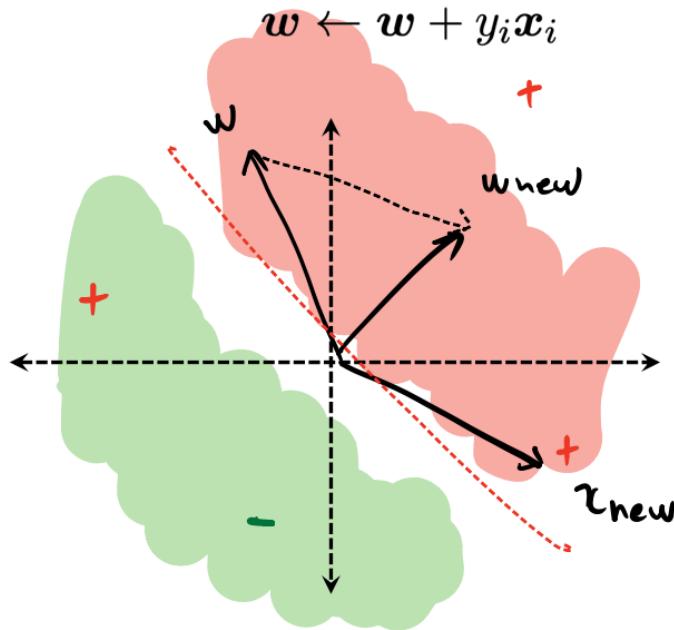
if  $x_i \neq 0$

$$y_i (w')^T x_i > y_i w^T x_i$$

# Perceptron algorithm: visually

Repeat:

- Pick a data point  $x_i$  uniformly at random
- If  $\text{sgn}(\mathbf{w}^T \mathbf{x}_i) \neq y_i$

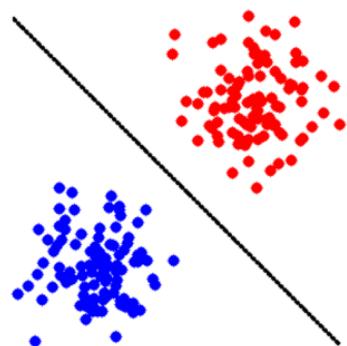


Related to question in class: if there are multiple ways to classify data:  
This might be better  
Perceptron itself would find any of these hyperpl

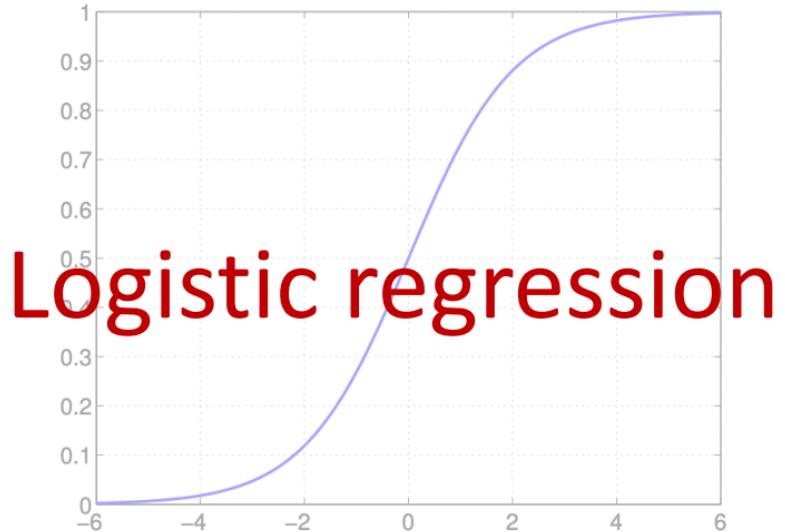
# HW1: Theory for perceptron!

(HW 1) If training set is linearly separable

- Perceptron *converges in a finite number of steps*
- training error is 0



There are also guarantees when the data are not linearly separable.



## Logistic loss

$$F(\omega) = \frac{1}{n} \sum_{i=1}^n \ell_{\log}(y_i \omega^\top x_i)$$
$$= \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y_i \omega^\top x_i})$$

# Predicting probabilities

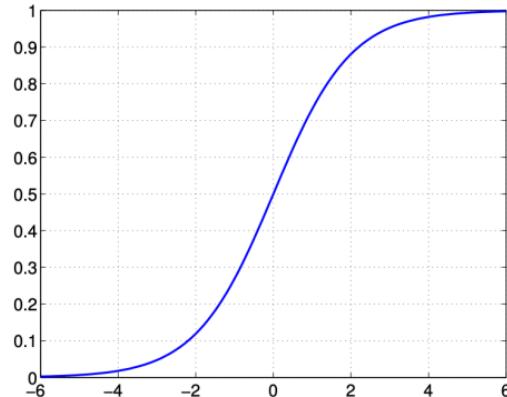
Instead of  $\{ \pm 1 \}$ 's, predict the probability.

(i.e. regression on probability)

Sigmoid + linear model:

$$\text{IP}(y=+1 | x, w) = \sigma(w^T x)$$

where  $\sigma(z) = \frac{1}{1+e^{-z}}$   
(sigmoid)



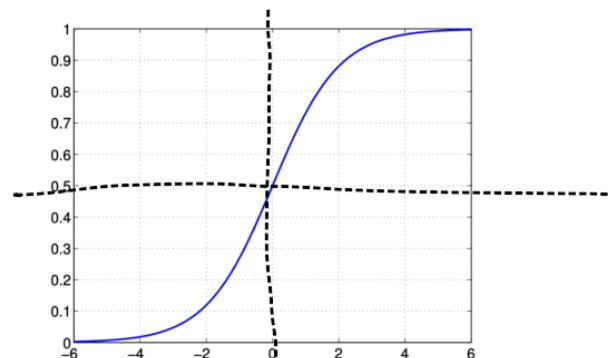
# The sigmoid function

**Properties** of sigmoid  $\sigma(z) = \frac{1}{1+e^{-z}}$

- between 0 and 1 (good as probability)
- $\sigma(\mathbf{w}^T \mathbf{x}) \geq 0.5 \Leftrightarrow \mathbf{w}^T \mathbf{x} \geq 0$ , consistent with predicting the label with  $\text{sgn}(\mathbf{w}^T \mathbf{x})$
- larger  $\mathbf{w}^T \mathbf{x} \Rightarrow$  larger  $\sigma(\mathbf{w}^T \mathbf{x}) \Rightarrow$  higher *confidence* in label 1
- $\sigma(z) + \sigma(-z) = 1$  for all  $z$
- Therefore, the probability of label  $-1$  is

$$\begin{aligned}\mathbb{P}(y = -1 \mid \mathbf{x}; \mathbf{w}) &= 1 - \mathbb{P}(y = +1 \mid \mathbf{x}; \mathbf{w}) \\ &= 1 - \sigma(\mathbf{w}^T \mathbf{x}) = \sigma(-\mathbf{w}^T \mathbf{x})\end{aligned}$$

Therefore, we can model  $\mathbb{P}(y \mid \mathbf{x}; \mathbf{w}) = \sigma(y \mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-y \mathbf{w}^T \mathbf{x}}}$



# Maximum likelihood estimation

*What we observe are labels, not probabilities.*

Take a **probabilistic view**

- assume data is independently generated in this way by some  $\mathbf{w}$
- perform Maximum Likelihood Estimation (MLE)

Specifically, what is the probability of seeing labels  $y_1, \dots, y_n$  given  $x_1, \dots, x_n$ , as a function of some  $\mathbf{w}$ ?

$$P(\mathbf{w}) = \prod_{i=1}^N \mathbb{P}(y_i \mid \mathbf{x}_i; \mathbf{w})$$

**MLE:** find  $\mathbf{w}^*$  that **maximizes the probability**  $P(\mathbf{w})$

## Maximum likelihood solution

$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} P(\mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^n \mathbb{P}(y_i | \mathbf{x}_i; \mathbf{w})$$

$$= \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^n \ln \mathbb{P}(y_i | \mathbf{x}_i; \mathbf{w})$$

$$= \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n -\ln \mathbb{P}(y_i | \mathbf{x}_i; \mathbf{w})$$

$$= \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n \ln(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})$$

$$= \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n \ell_{\text{logistic}}(y_i \mathbf{w}^T \mathbf{x}_i)$$

$$= \operatorname{argmin}_{\mathbf{w}} F(\mathbf{w})$$

Minimizing logistic loss is exactly doing MLE for the sigmoid model!

## SGD to logistic loss

$$\begin{aligned}
 \mathbf{w} &\leftarrow \mathbf{w} - \eta \tilde{\nabla} F(\mathbf{w}) \\
 &= \mathbf{w} - \eta \nabla_{\mathbf{w}} \ell_{\text{logistic}}(y_i \mathbf{w}^T \mathbf{x}_i) \\
 &= \mathbf{w} - \eta \left( \frac{\partial \ell_{\text{logistic}}(z)}{\partial z} \Big|_{z=y_i \mathbf{w}^T \mathbf{x}_i} \right) y_i \mathbf{x}_i \\
 &= \mathbf{w} - \eta \left( \frac{-e^{-z}}{1 + e^{-z}} \Big|_{z=y_i \mathbf{w}^T \mathbf{x}_i} \right) y_i \mathbf{x}_i \\
 &= \mathbf{w} + \eta \sigma(-y_i \mathbf{w}^T \mathbf{x}_i) y_i \mathbf{x}_i \\
 &= \mathbf{w} + \eta \mathbb{P}(-y_i \mid \mathbf{x}_i; \mathbf{w}) y_i \mathbf{x}_i
 \end{aligned}$$

$$\mathbb{E}[\tilde{\nabla} F(\mathbf{w})] = \nabla F(\mathbf{w})$$

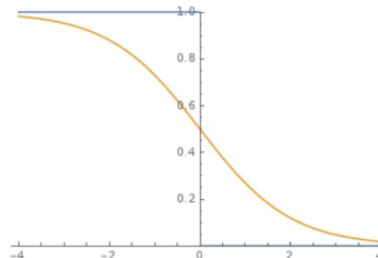
( $i$  is drawn uniformly from  $[n]$ )  
(Chain rule)

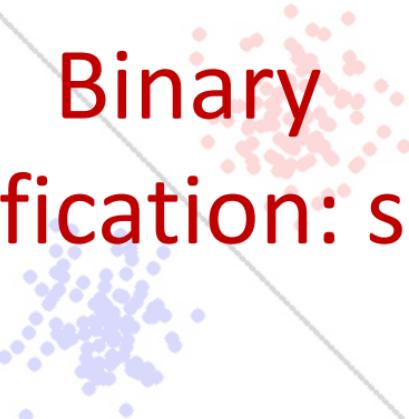
$$\frac{\partial \log(1 + e^{-z})}{\partial z} = \frac{-e^{-z}}{1 + e^{-z}}$$

$$\sigma(-z) = 1 - \sigma(z) = 1 - \frac{1}{1 + e^{-z}} = \frac{e^{-z}}{1 + e^{-z}}$$

This is a *soft version of Perceptron!*

$\mathbb{P}(-y_i \mid \mathbf{x}_i; \mathbf{w})$  versus  $\mathbb{I}[y_i \neq \text{sgn}(\mathbf{w}^T \mathbf{x}_i)]$





Binary  
classification: so far

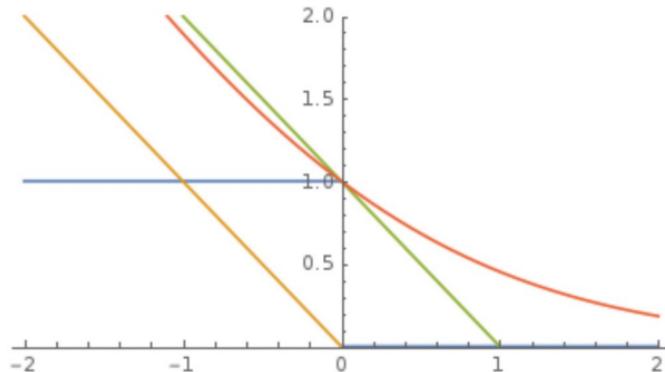
# Supervised learning in one slide

- Loss function:** What is the right loss function for the task?
- Representation:** What class of functions should we use?
- Optimization:** How can we efficiently solve the empirical risk minimization problem?
- Generalization:** Will the predictions of our model transfer gracefully to unseen examples?

*All related! And the fuel which powers everything is **data**.*

## Loss function

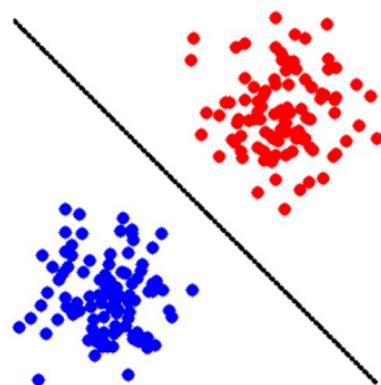
Use a **convex surrogate loss**



- **perceptron loss**  $\ell_{\text{perceptron}}(z) = \max\{0, -z\}$  (used in Perceptron)
- **hinge loss**  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$  (used in SVM and many others)
- **logistic loss**  $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z))$  (used in logistic regression; the base of log doesn't matter)

## Representation

Definition: The **function class of separating hyperplanes** is defined as  
 $\mathcal{F} = \{f(x) = \text{sign}(\mathbf{w}^T \mathbf{x}): \mathbf{w} \in \mathbb{R}^d\}$ .



## Optimization

Empirical risk minimization (ERM) problem:

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(y_i \mathbf{w}^T \mathbf{x}_i)$$

Solve using a suitable optimization algorithm:

- **GD:**  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla F(\mathbf{w})$
- **SGD:**  $\mathbf{w} \leftarrow \mathbf{w} - \eta \tilde{\nabla} F(\mathbf{w})$   $(\mathbb{E}[\tilde{\nabla} F(\mathbf{w})] = \nabla F(\mathbf{w}))$
- **Newton:**  $\mathbf{w} \leftarrow \mathbf{w} - (\nabla^2 F(\mathbf{w}))^{-1} \nabla F(\mathbf{w})$

## Generalization

**Rich theory!** Let's see a glimpse 😊



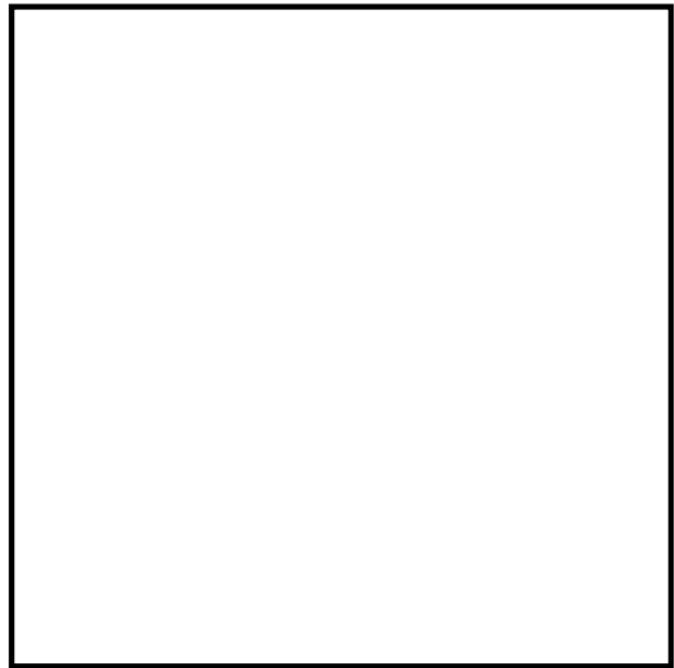
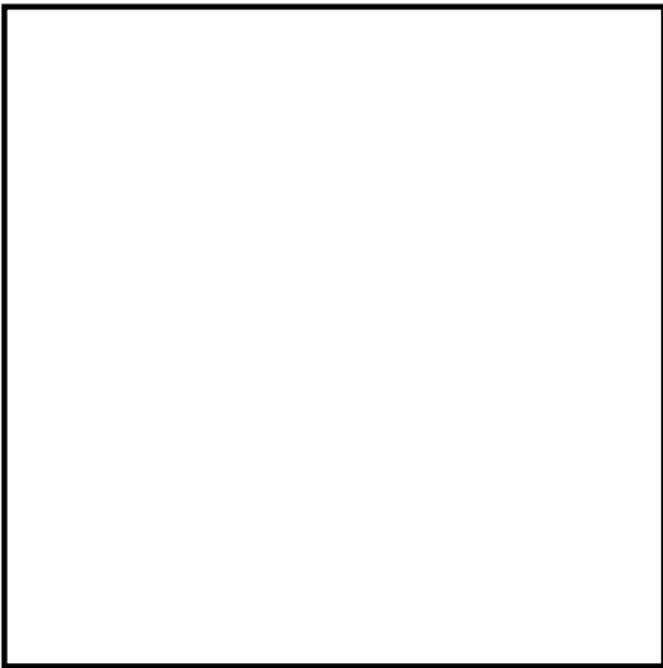
# Generalization

## Reviewing definitions

- Input space:  $\mathcal{X}$
- Output space:  $\mathcal{Y}$
- Predictor:  $f(\mathbf{x}) : \mathcal{X} \rightarrow \mathcal{Y}$
- Distribution  $D$  over  $(\mathbf{x}, y)$ .
- Let  $D^n$  denote the distribution of  $n$  samples  $\{(\mathbf{x}_i, y_i), i \in [n]\}$  drawn i.i.d. from  $D$ .
- Risk of a predictor  $f(\mathbf{x})$  is  $R(f) = \mathbb{E}_{(\mathbf{x}, y) \sim D} [\ell(f(\mathbf{x}), y)]$
- Consider the 0-1 loss,  $\ell(f(\mathbf{x}, y)) = \mathbb{1}(f(\mathbf{x}) \neq y)$ .

**Next time, we'll see some  
generalization theory!**

## Intuition: When does ERM generalize?













## Relaxing our assumptions

- We assumed that the function class is finite-sized. Results can be extended to **infinite function classes** (such as separating hyperplanes).
- We considered 0-1 loss. Can extend to **real-valued loss** (such as for regression).
- We assumed realizability. Can prove similar theorem which guarantees small generalization gap **without realizability** (but with an  $\epsilon^2$  instead of  $\epsilon$  in the denominator). This is called agnostic learning.

## Rule of thumb for generalization

Suppose the functions  $f$  in our function class  $\mathcal{F}$  have  $d$  parameters which can be set. Assume we discretize these parameters so they can take  $W$  possible values each. How much data do we need to have small generalization gap?

A useful rule of thumb: to guarantee generalization, make sure that your training data set size  $n$  is at least linear in the number  $d$  of free parameters in the function that you're trying to learn.