



Some Basics of Machine Learning

Hongbo Li, Postdoc Scholar

Time: 2:00pm – 3:00pm, June 10, 2025

Slide Credits

1. *Introduction to Machine Learning*, 10-401, by Maria-Florina Balcan, Carnegie Mellon University:
<https://www.cs.cmu.edu/~ninamf/courses/401sp18/lectures.shtml>
2. *Machine Learning*, CS4/5780, by Kilian Weinberger, Cornell University: <https://www.cs.cornell.edu/courses/cs4780/2017sp/>
3. *Intro to Machine Learning*, CS4/CS5780, by Wen Sun, Cornell University: <https://www.cs.cornell.edu/courses/cs4780/2023fa/>
4. *Deep Learning for Computer Vision*, CS231n, by Fei-Fei Li et. al., Stanford University: <https://cs231n.stanford.edu>

Index

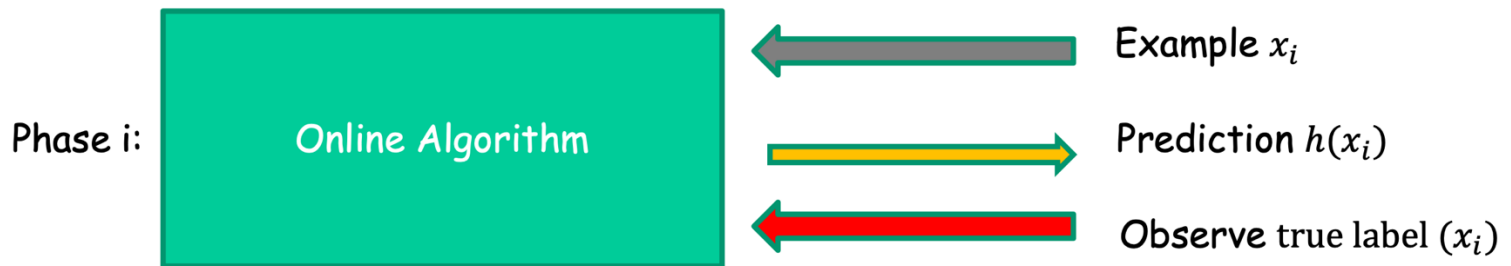
- Part I: Linear Classifier
- Part II: Gradient Descent (and Beyond)
- Part III: Linear Regression
- Part IV: Convolutional Neural Networks (CNNs)

Part I: Linear Classifier

Online Learning Model

- Examples arrive **sequentially**.

For $i=1, 2, \dots, :$

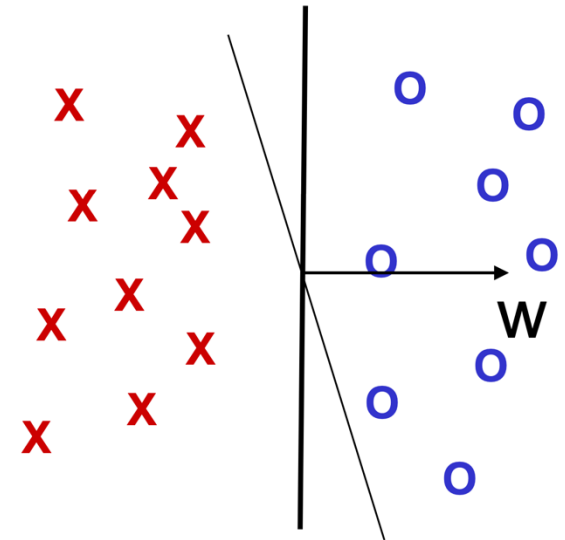


- Objective: **Minimize** the number of **mistakes**.
- Applications:
 - Email classification (distribution of both spam and regular mail changes over time).
 - Recommendation systems.
 - ...

Figure from Balcan (2018)

Linear Separators

- Feature space $X \in \mathbb{R}^d$.
- Hypothesis class of linear decision surfaces in \mathbb{R}^d .
 - $h(x) = \mathbf{W}^T x + w_0$, where $\mathbf{W} = (w_1, \dots, w_d) \in \mathbb{R}^d$
 - If $h(x) \geq 0$, then label x as $+$, otherwise label it as $-$



Trick: Without loss of generality $w_0 = 0$.

Proof: Can simulate a non-zero threshold with a dummy input feature x_0 that is always set up to 1.

- $x = (x_1, \dots, x_d) \rightarrow \tilde{x} = (x_1, \dots, x_d, 1)$
- $\mathbf{W}^T x + w_0 \geq 0$ iff $\widetilde{\mathbf{W}}^T \tilde{x} \geq 0$, where $\widetilde{\mathbf{W}} = (w_1, \dots, w_d, w_0)$

Figure from Balcan (2018)

Linear Separators: **Perceptron Algorithm**

- Set $t = 1$, start with the all zero vector $\mathbf{W}_1 = (0, \dots, 0)$.
- Given example \mathbf{x} , predict positive iff $\mathbf{W}_t^T \mathbf{x} \geq 0$.
- On a mistake, update as follows:
 - Mistake on positive, then update $\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t + \mathbf{x}$
 - Mistake on negative, then update $\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t - \mathbf{x}$

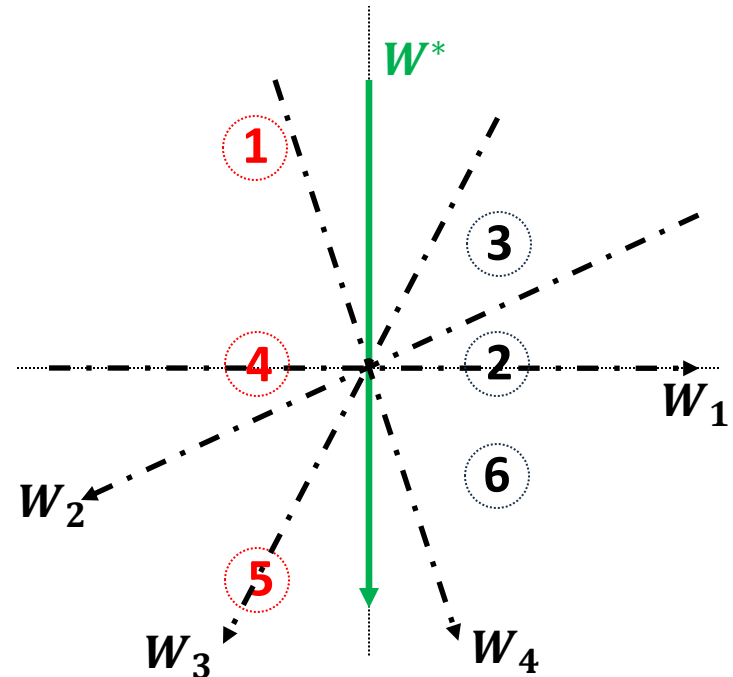
Natural **greedy** procedure:

- If true label of \mathbf{x} is +1 and \mathbf{W}_t **incorrect** on \mathbf{x} , we have $\mathbf{W}_t^T \mathbf{x} < 0$.
- By Perceptron Algorithm, we have $\mathbf{W}_{t+1}^T \mathbf{x} \leftarrow \mathbf{W}_t^T \mathbf{x} + \mathbf{x}^T \mathbf{x} = \mathbf{W}_t^T \mathbf{x} + \|\mathbf{x}\|^2$.
- Then, there will be more chance that \mathbf{W}_{t+1} classifies \mathbf{x} correctly.
- Similar for mistakes on negative examples.

Perceptron Algorithm: Practice

Example: $W^* = (1, 0)$

- | | | | | |
|----|------------|---|-----------------|---|
| 1. | $(-1, 2)$ | - | $W_1 = (0, 0)$ | ✗ |
| 2. | $(1, 0)$ | + | $W_2 = (1, -2)$ | ✓ |
| 3. | $(1, 1)$ | + | $W_2 = (1, -2)$ | ✗ |
| 4. | $(-1, 0)$ | - | $W_3 = (2, -1)$ | ✓ |
| 5. | $(-1, -2)$ | - | $W_3 = (2, -1)$ | ✗ |
| 6. | $(1, -1)$ | + | $W_4 = (3, 1)$ | ✓ |



Algorithm:

- Set $t = 1$, start with the all zero vector $W_1 = (0, \dots, 0)$.
- Given example x , predict positive iff $W_t^T x \geq 0$.
- On a mistake, update as follows:
 - Mistake on positive, then update $W_{t+1} \leftarrow W_t + x$
 - Mistake on negative, then update $W_{t+1} \leftarrow W_t - x$

$$W_1 = (0, 0)$$

$$W_2 = W_1 - (-1, 2) = (1, -2)$$

$$W_3 = W_2 + (1, 1) = (2, -1)$$

$$W_4 = W_3 - (-1, -2) = (3, 1)$$

Geometric Margin

Definition: The **margin** of **example x** w.r.t a linear separator **W** is the distance from **x** to the plane **$W^T x = 0$** (or the negative if on wrong side)

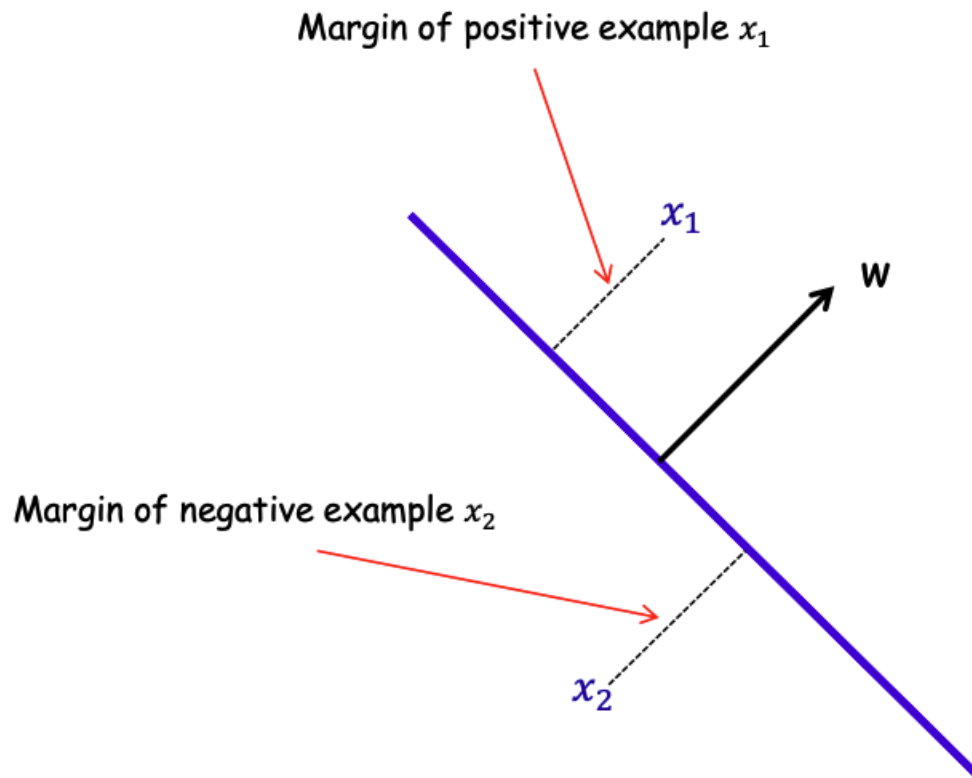


Figure from Balcan (2018)

Geometric Margin

Definition: The **margin** of **example x** w.r.t a linear separator W is the distance from x to the plane $W^T x = 0$ (or the negative if on wrong side)

Definition: The **margin γ_W** of a set of example S w.r.t a linear separator W is the smallest margin over points $x \in S$.

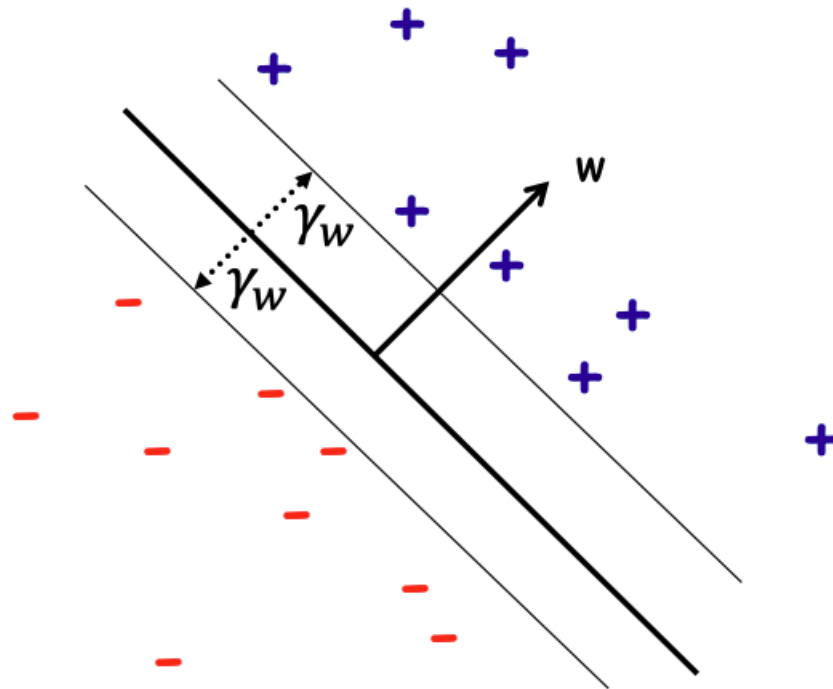


Figure from Balcan (2018)

Geometric Margin

Definition: The **margin** of **example x** w.r.t a linear separator W is the distance from x to the plane $W^T x = 0$ (or the negative if on wrong side)

Definition: The **margin γ_W** of a set of example S w.r.t a linear separator W is the smallest margin over points $x \in S$.

Definition: The **margin γ** of a set of example S is the **maximum γ_W** over all linear separators w .

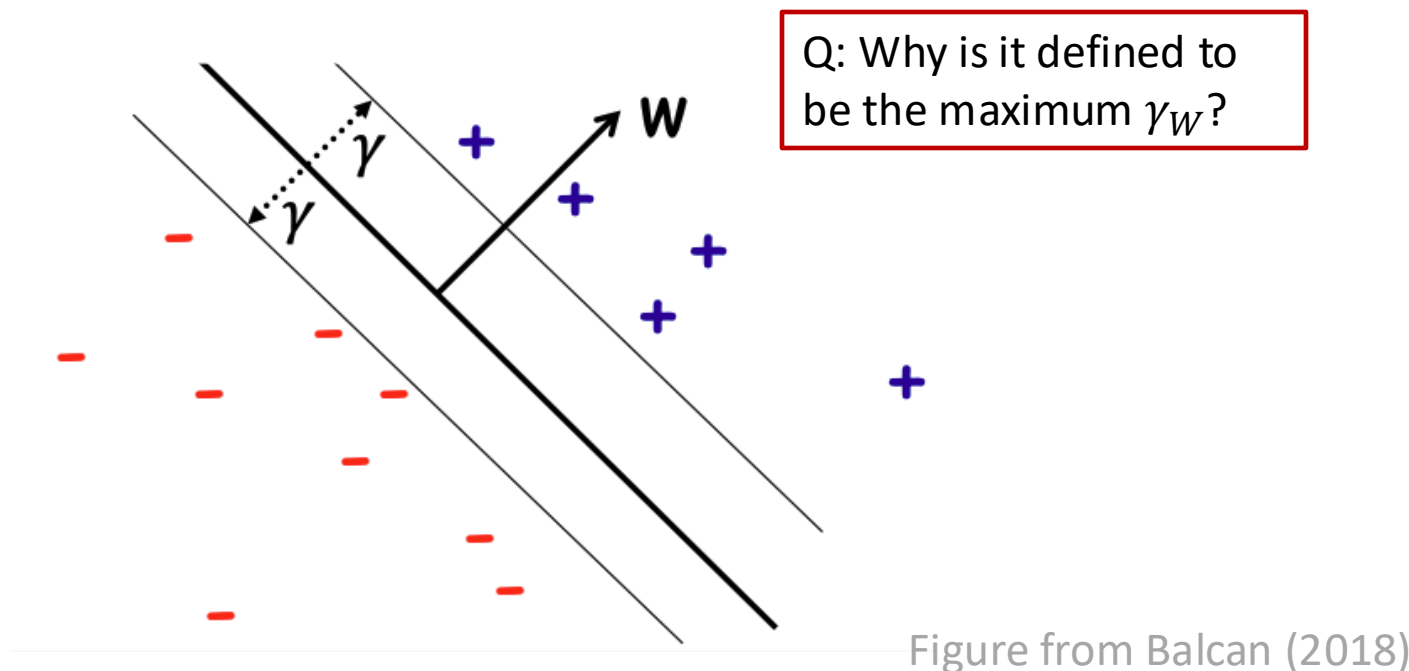


Figure from Balcan (2018)

Perceptron: Mistake Bound

Guarantee: If data has margin γ and all points inside a ball of radius R , then Perceptron makes $\leq (R/\gamma)^2$ mistakes.

(Normalized margin: multiplying all points by 100, or dividing all points by 100, doesn't change the number of mistakes: algo is invariant to scaling.)

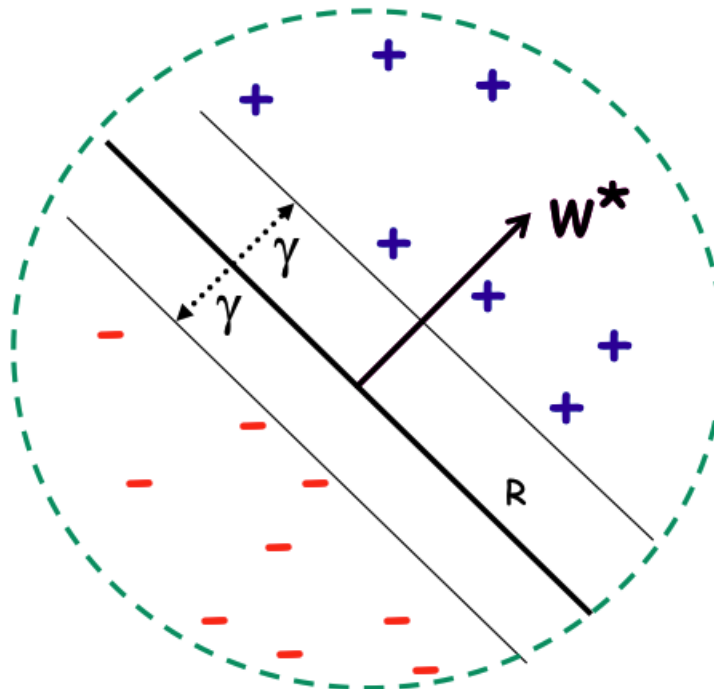


Figure from Balcan (2018)

Perceptron: Mistake Bound

Guarantee: If data has margin γ and all points inside a ball of radius R , then Perceptron makes $\leq (R/\gamma)^2$ mistakes.

Update rule:

- Mistake on positive: $\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t + x$
- Mistake on negative: $\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t - x$

Proof:

Idea: analyze $\mathbf{W}_t^T \mathbf{W}^*$ and $\|\mathbf{W}_t\|$, where \mathbf{W}^* is the max-margin separator with $\|\mathbf{W}^*\| = 1$.

- Claim 1: $\mathbf{W}_{t+1}^T \mathbf{W}^* \geq \mathbf{W}_t^T \mathbf{W}^* + \gamma$. (because $\mathbf{W}_t^T x \geq \gamma$)
- Claim 2: $\|\mathbf{W}_{t+1}\|^2 \leq \|\mathbf{W}_t\|^2 + R^2$. (by Pythagorean Theorem)

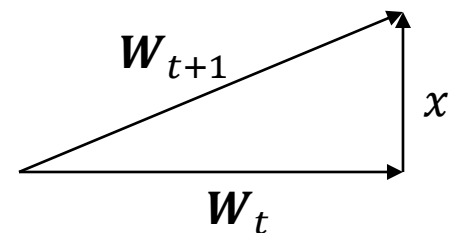
After M mistakes:

$$\mathbf{W}_{t+1}^T \mathbf{W}^* \geq \gamma M \quad (\text{by Claim 1})$$

$$\|\mathbf{W}_{M+1}\| \leq R\sqrt{M} \quad (\text{by Claim 2})$$

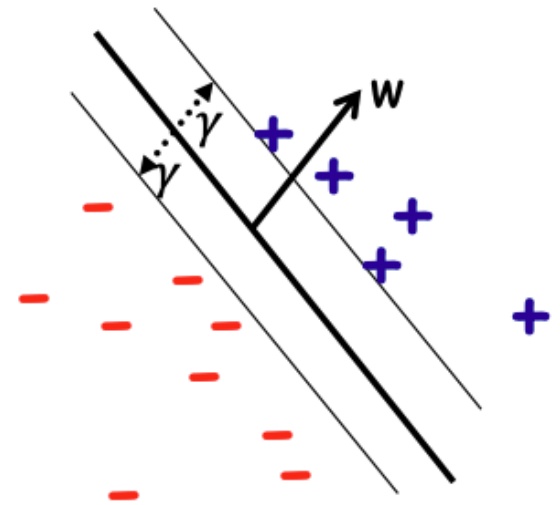
$$\mathbf{W}_{t+1}^T \mathbf{W}^* \leq \|\mathbf{W}_{M+1}\| \quad (\text{since } \mathbf{W}^* \text{ is unit length})$$

In summary, $\gamma M \leq \mathbf{W}_{t+1}^T \mathbf{W}^* \leq \|\mathbf{W}_{M+1}\| \leq R\sqrt{M}$. So, $M \leq (R/\gamma)^2$.



Margin Important Theme in ML

- If large margin, the number of mistakes Perceptron makes is small (independent on the dimension of the ambient space)
- Large margin can help prevent **overfitting**.
- **Support Vector Machines (SVMs)** directly optimize for the maximum margin separator:
 - Input: $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$;
 - Find: some \mathbf{W} and maximum γ where:
 - $\|\mathbf{W}\| = 1$
 - For all $i \in \{1, \dots, m\}, y_i \mathbf{W}^T \geq \gamma$
 - Output: maximum margin separator over S .



Part II: Gradient Descent (and Beyond)

Taylor Expansion

- How can you minimize a function $l(\mathbf{W})$ if you don't know much about it?
- The trick is to **assume it is much simpler** than it really is.
- This can be done with **Taylor's approximation**.
- Given a small norm $\|s\|_2$ (i.e., $\mathbf{W} + s$ is very close to \mathbf{W}), we can approximate the function $l(\mathbf{W} + s)$ by its first and second derivatives:

$$l(\mathbf{W} + s) \approx l(\mathbf{W}) + g(\mathbf{W}) \cdot s \text{ (Gradient Descent)}$$

$$l(\mathbf{W} + s) \approx l(\mathbf{W}) + g(\mathbf{W}) \cdot s + \frac{1}{2} s^T \cdot H(\mathbf{W}) \cdot s \text{ (Newton's Method)}$$

Here $g(x) = \nabla l(\mathbf{W})$ is the gradient and $H(x) = \nabla^2 l(\mathbf{W})$ is the Hessian of l .

Gradient Descent (GD)

- In GD, we only use the gradient (**first order**).
- We assume the function l around \mathbf{W} is **linear** and behaves like $l(\mathbf{W}) + g(\mathbf{W}) \cdot s$.
- Our objective is to find a vector s that **minimizes** function l .
- In steepest descent we simply set

$$s = -\eta \cdot g(\mathbf{W})$$

for some small $\eta > 0$.

- It is straight-forward to prove that in this case $l(\mathbf{W} + s) < l(\mathbf{W})$:

$$\begin{aligned} & l(\mathbf{W} + (-\eta \cdot g(\mathbf{W}))) \\ & \approx l(\mathbf{W}) - \eta \cdot g(\mathbf{W})^T g(\mathbf{W}) \\ & < l(\mathbf{W}) \end{aligned}$$

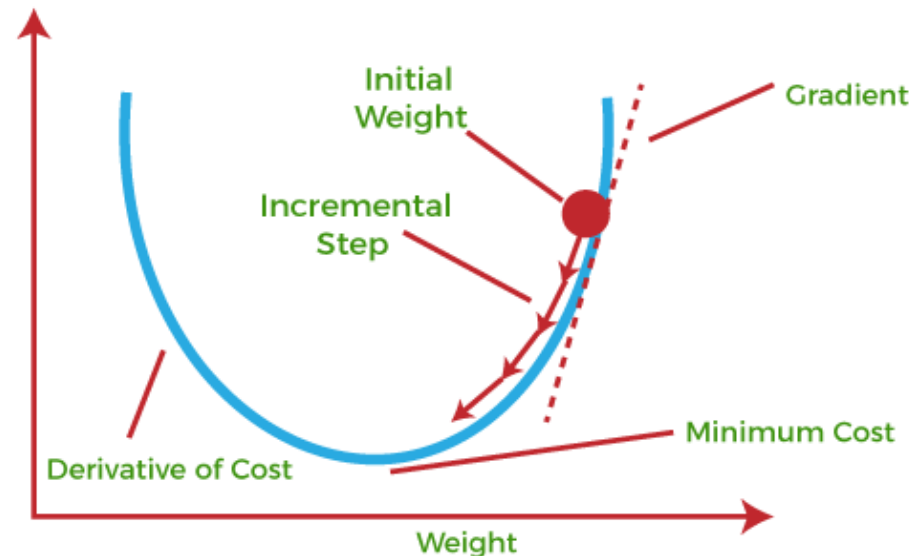
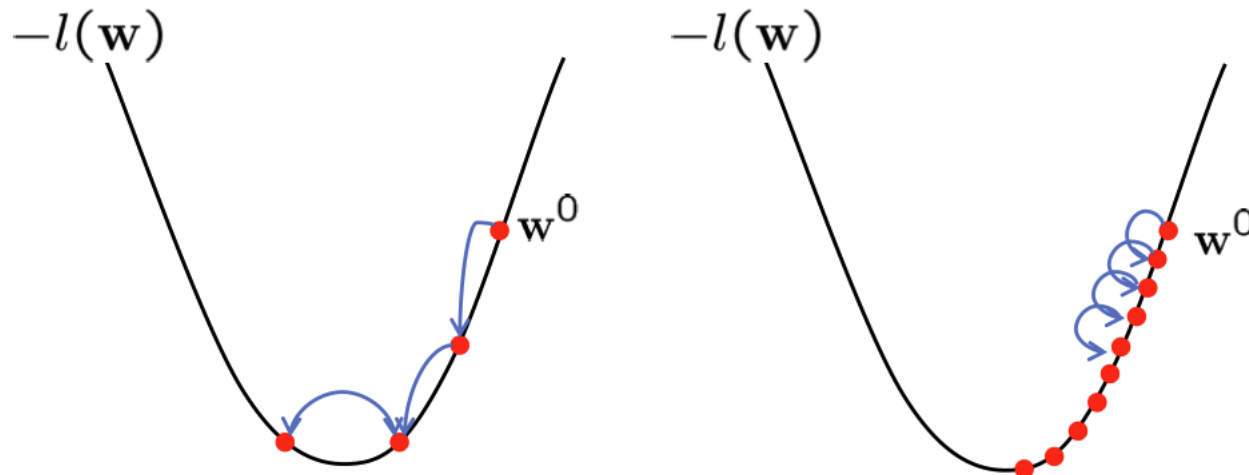


Figure from Kharkar (2023)

GD: Learning Rate

- Setting the **learning rate** $\eta > 0$ is a dark art.
 - Large η** \Rightarrow Fast convergence but larger residual error $\|W^{t+1} - W^t\|_2$, with possible oscillations.
 - Small η** \Rightarrow Slow convergence but small residual error.



- A safe (but sometimes slow) choice is to set $\eta = \frac{1}{t}$, which guarantees that it will eventually become small enough to converge.

Figure from Balcan (2018)

GD: In Practice

In ML, the loss we minimize typically has some special form, e.g.,

$$l(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i(\mathbf{W}^T \mathbf{x}_i)))$$

Average over n data points

To compute the gradient $\nabla l(\mathbf{W})$, we need to enumerate all n training data points, which **can be very slow!**

Stochastic GD (**SGD**) to Rescue

In ML, the loss we minimize typically has some special form, e.g.,

$$l(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i(\mathbf{W}^T \mathbf{x}_i)))$$

Average over n data points

Idea: Randomly sample a data point (x, y) , use $\nabla l(x, y; \mathbf{W})$ to replace $\nabla l(\mathbf{W})$.

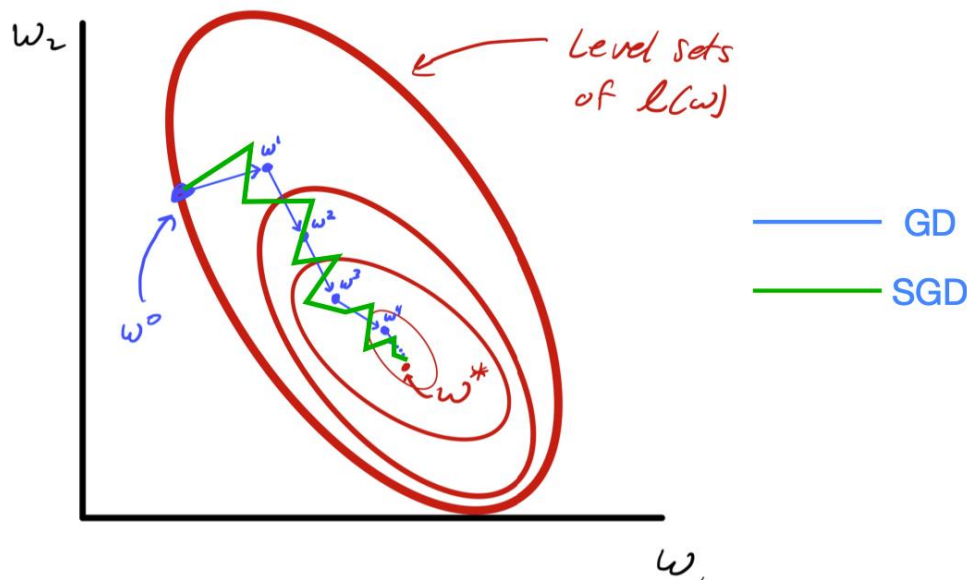
Stochastic GD (SGD)

- Goal: minimize $l(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n l(x_i, y_i; \mathbf{W})$
- Initialize: $\mathbf{W}^0 \in \mathbb{R}^d$ randomly
- Iterate until convergence:
 1. Randomly sample a point (x_i, y_i) from the n data points
 2. Compute noisy gradient $\tilde{g}^t = \nabla l(x_i, y_i; \mathbf{W})|_{\mathbf{W}=\mathbf{W}^t}$
 3. Update (GD): $\mathbf{W}^{t+1} = \mathbf{W}^t - \eta \tilde{g}^t$

Why Can SGD Work?

Claim: the random noisy gradient is an **unbiased estimate** of the true gradient

$$\mathbb{E}[\nabla l(x_i, y_i; \mathbf{W})] = \frac{1}{n} \sum_{i=1}^n \nabla l(x_i, y_i; \mathbf{W}) = \nabla \left[\frac{1}{n} \sum_{i=1}^n l(x_i, y_i; \mathbf{W}) \right] = \nabla l(\mathbf{W})$$



Q: Which one is faster?

A: SGD is **slower in theory**, but is often **much faster in practice**.

Figure from Sun (2022)

Part III: Linear Regression

From Discrete to Continuous Labels

Classification:



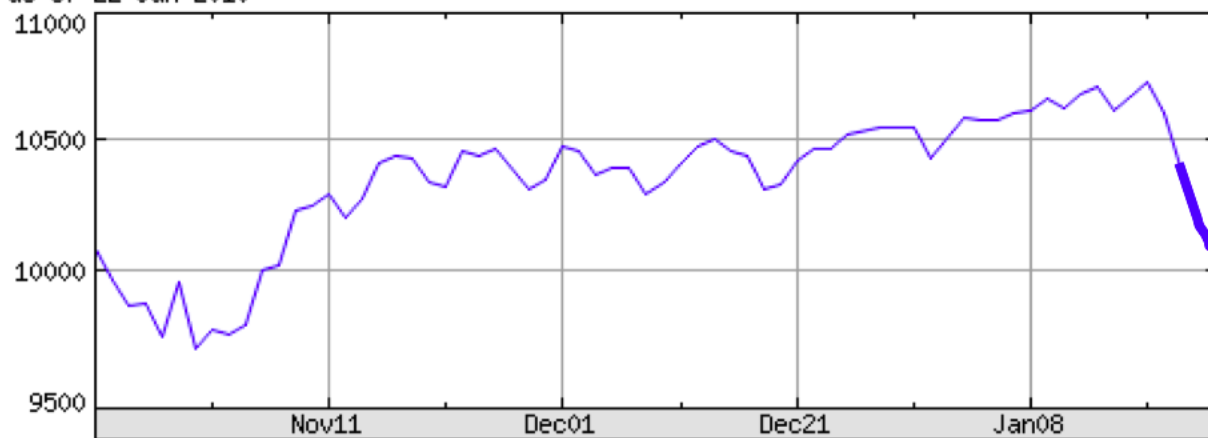
Sports
Science
News

$X = \text{Document}$

$Y = \text{Topic}$

Regression:

DJ INDU AVERAGE (DOW JONES & CO)
as of 22-Jan-2010



$Y = ?$

$X = \text{Feb01}$

Copyright 2010 Yahoo! Inc.

<http://finance.yahoo.com/>

Figure from Balcan (2018)

Supervised Learning

Goal: Construct a predictor $f: X \rightarrow Y$ to minimize a risk (error measure) $\text{err}(f)$.



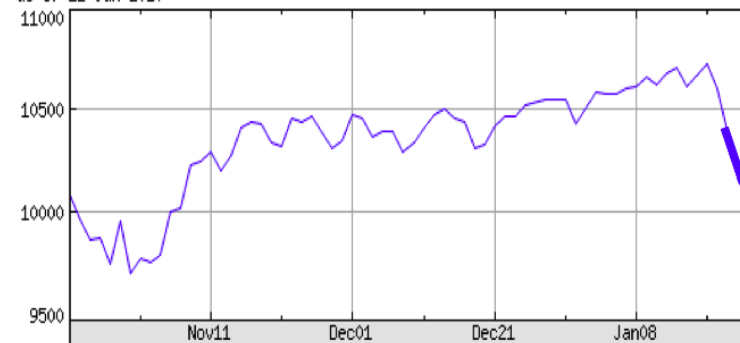
X = Document



**Sports
Science
News**

Y = Topic

DJ INDU AVERAGE (DOW JONES & CO)
as of 22-Jan-2010



Copyright 2010 Yahoo! Inc.

<http://finance.yahoo.com/>

X = Feb01

Classification:

$$\text{err}(f) = P(f(X) \neq Y)$$

Probability of error

Regression:

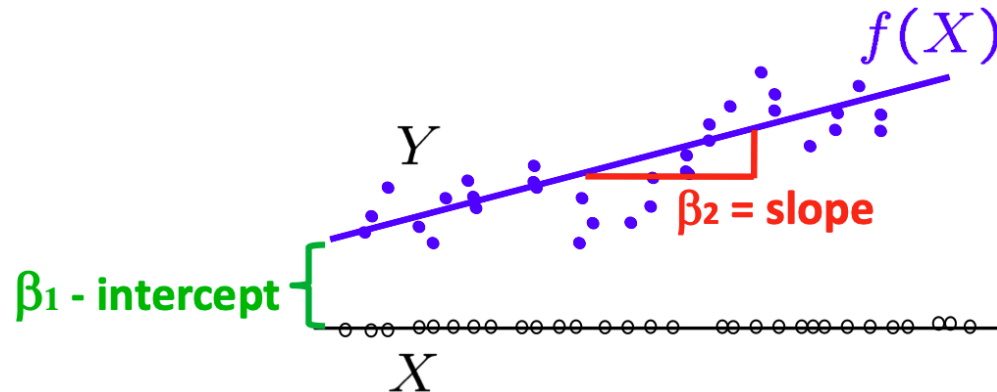
$$\text{err}(f) = E[(f(X) - Y)^2]$$

Mean Squared Error

Figure from Balcan (2018)

Linear Regression

- Unit-variate case: $f(X) = \beta_1 + \beta_2 X$.



- Multi-variate case:

$$f(X) = f(X^{(1)}, \dots, X^{(p)}) = \beta_1 X^{(1)} + \beta_2 X^{(2)} + \dots + \beta_p X^{(p)}$$

$$= X\beta, \text{ where } X = [X^{(1)} \dots X^{(p)}], \beta = [\beta_1 \dots \beta_p]^T$$
- Least square estimator: $\hat{f} = \arg \min_{f \in F} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2$, where F is the class of linear functions.

Figure from Balcan (2018)

Least Squares Estimator

- $\hat{f} = \arg \min_{f \in F} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2$



- $\hat{\beta} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n (X_i \beta - Y_i)^2$
 $= \arg \min_{\beta} \frac{1}{n} (\mathbf{X}\beta - \mathbf{Y})^T (\mathbf{X}\beta - \mathbf{Y})$

$$\mathbf{X} = \begin{bmatrix} X_1 \\ \dots \\ X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \dots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \dots & X_n^{(p)} \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} Y_1 \\ \dots \\ Y_n \end{bmatrix}.$$

Least Squares Estimator

- $\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{X}\beta - \mathbf{Y})^T (\mathbf{X}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$
- $J(\beta) = (\mathbf{X}\beta - \mathbf{Y})^T (\mathbf{X}\beta - \mathbf{Y})$
- $\frac{\partial J(\beta)}{\partial \beta} \big|_{\hat{\beta}} = 0$



- $(\mathbf{X}^T \mathbf{X}) \hat{\beta} = \mathbf{X}^T \mathbf{Y}$
- If $\mathbf{X}^T \mathbf{X}$ is invertible,

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad \hat{f}(\mathbf{X}) = \mathbf{X} \hat{\beta}$$

Least Squares Estimator: Verification

- $\hat{f}(\mathbf{X}) = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}\mathbf{Y}^T$
- We calculate

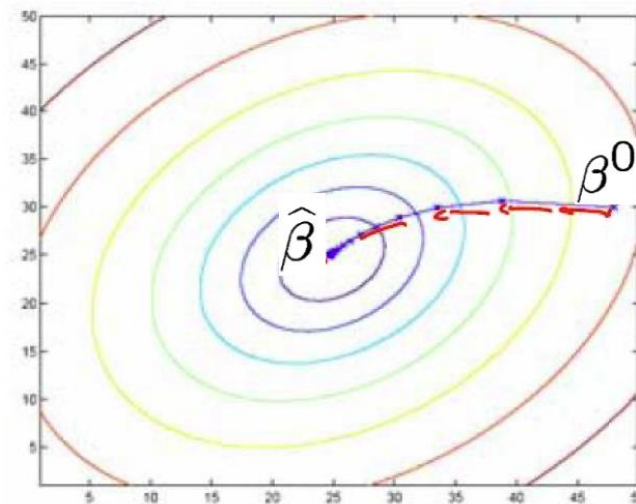
$$\mathbf{X}^T(\hat{f}(\mathbf{X}) - \mathbf{Y}) = \mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}\mathbf{Y}^T - \mathbf{X}\mathbf{Y}^T = \mathbf{0}.$$

Revisiting Gradient Descent

- Even when $\mathbf{X}^T \mathbf{X}$ is invertible, might be computationally expensive if \mathbf{X} is huge.

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{X}\beta - \mathbf{Y})^T (\mathbf{X}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

- Gradient Descent since $J(\beta)$ is convex
 - Initialize: β^0
 - Update: $\beta^{t+1} = \beta^t - \frac{\eta}{2} \mathbf{X}^T \frac{\partial J(\beta)}{\partial \beta} \Big|_t$
 $= \beta^t - \eta \mathbf{X}^T (\mathbf{X}\beta^t - \mathbf{Y})$
 - Stop: when some criterion met, e.g., fixed # iterations, or $\frac{\partial J(\beta)}{\partial \beta} \Big|_t < \varepsilon$.



Q: What about Stochastic GD for linear regression?

Figure from Balcan (2018)

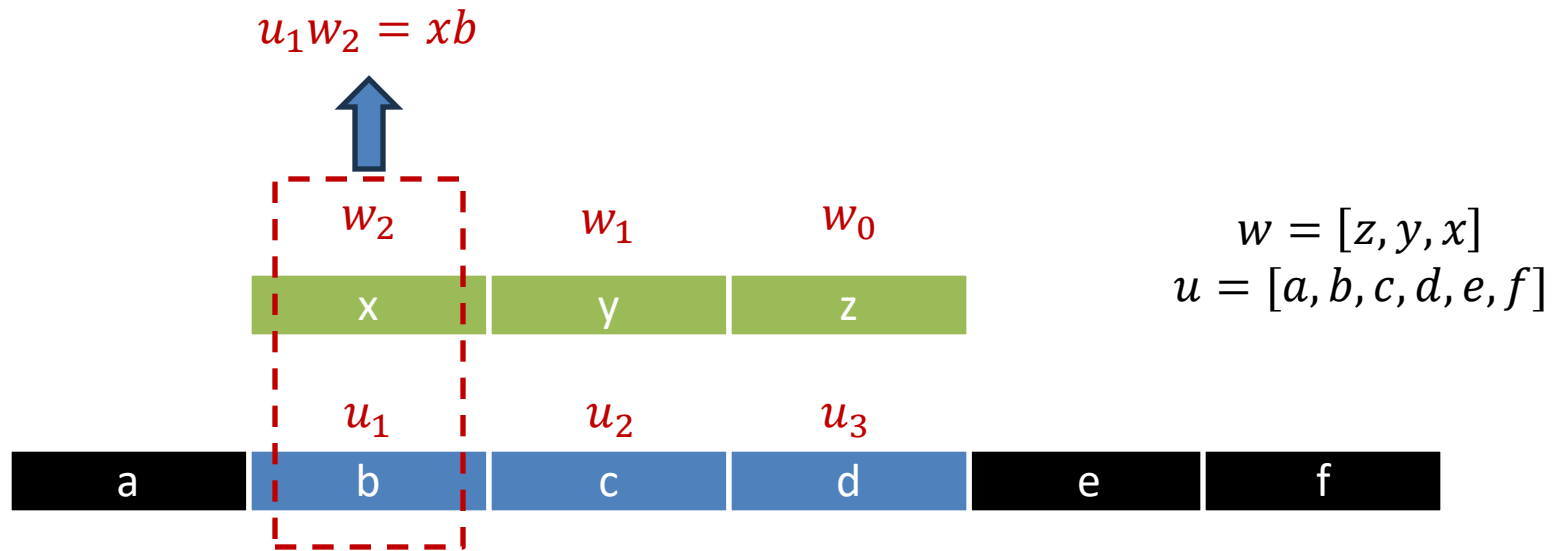
Part IV: Convolutional Neural Networks (CNNs)

Motivation of Convolution

- Suppose we track the location of a spaceship with a laser sensor. The laser sensor provides a single output $u(t)$, which is the position of the spaceship at second t .
- Suppose sensor is noisy. To obtain a less noisy estimate of the spaceship's position, we **average several measurements**. More recent measurements are more relevant, so we use a **weighted average** that gives more weight to recent measurements.
- Use a weighting function $w(a)$, where a is the age of a measurement. If we apply such a weighted average operation at every moment, we obtain a new function s providing a smoothed estimate of the position of the spaceship:

$$s_t = \sum_{a=-\infty}^{+\infty} u_a w_{t-a}$$

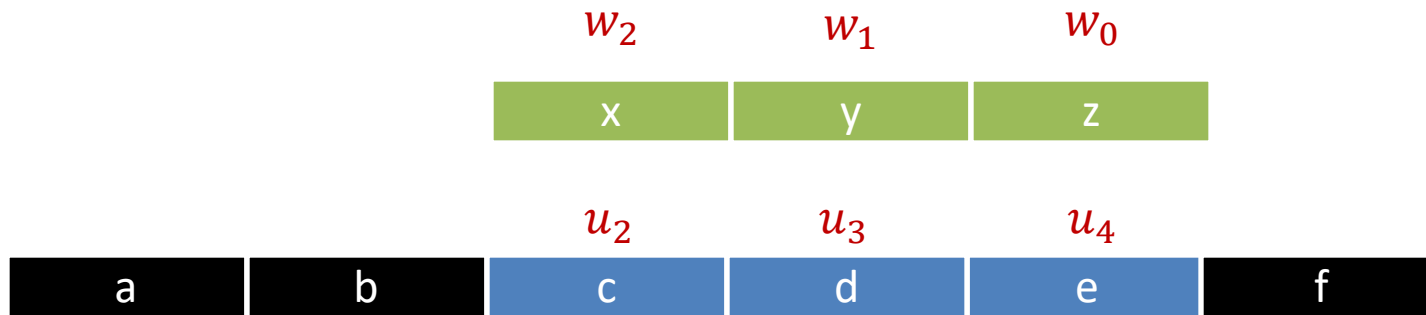
Illustration 1



$$s_t = \sum_{a=-\infty}^{+\infty} u_a w_{t-a} \quad \Rightarrow \quad s_3 = \sum_{a=1}^3 u_a w_{3-a} = u_1 w_2 + u_2 w_1 + u_3 w_0 = xb + yc + zd$$

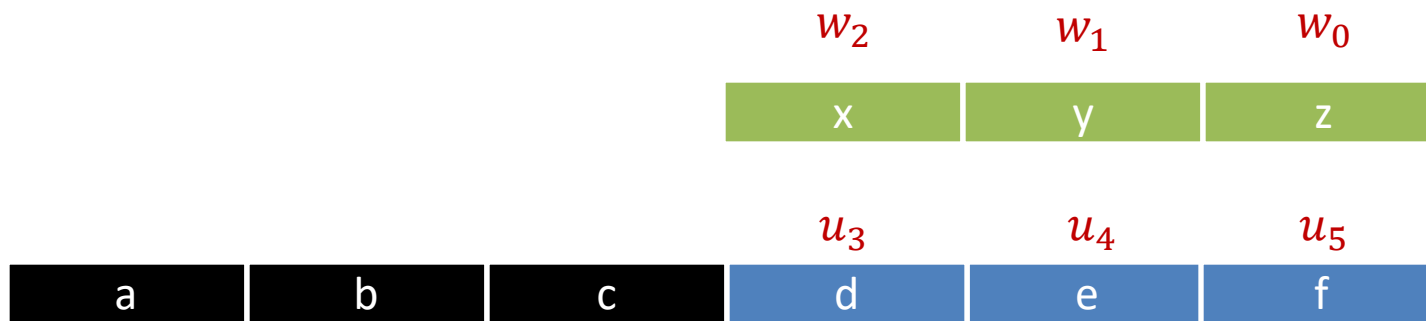
Q: What is s_3 here?

Illustration 1



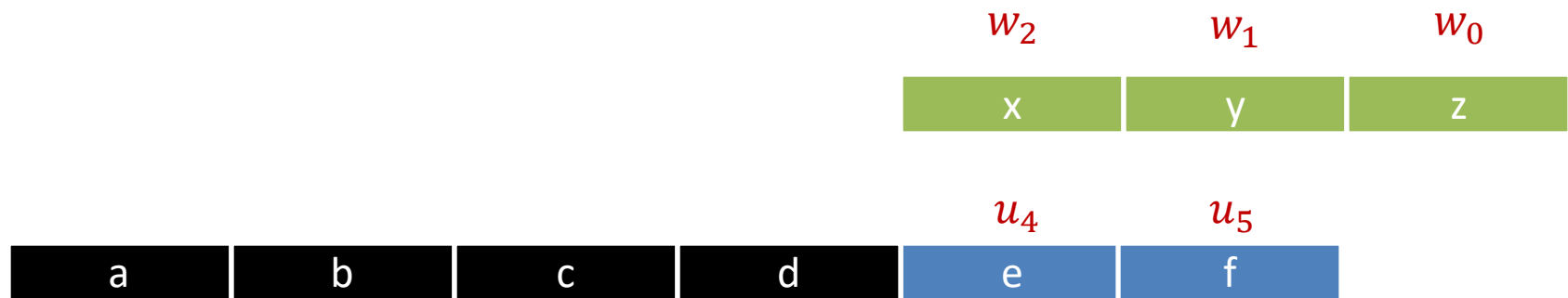
$$s_t = \sum_{a=-\infty}^{+\infty} u_a w_{t-a} \quad \Rightarrow \quad s_4 = \sum_{a=2}^4 u_a w_{4-a} = u_2 w_2 + u_3 w_1 + u_4 w_0 = xc + yd + ze$$

Illustration 1



$$s_t = \sum_{a=-\infty}^{+\infty} u_a w_{t-a} \quad \Rightarrow \quad s_5 = \sum_{a=3}^5 u_a w_{5-a} = u_3 w_2 + u_4 w_1 + u_5 w_0 = xd + ye + zf$$

Illustration 1



$$s_t = \sum_{a=-\infty}^{+\infty} u_a w_{t-a} \quad \Rightarrow \quad s_6 = \sum_{a=4}^5 u_a w_{6-a} = u_4 w_2 + u_5 w_1 = xe + yf$$

Illustration 1 as Matrix Multiplication

y	z					a
x	y	z				b
	x	y	z			c
		x	y	z		d
			x	y	z	e
				x	y	f

Illustration 2: Two-Dimensional Case

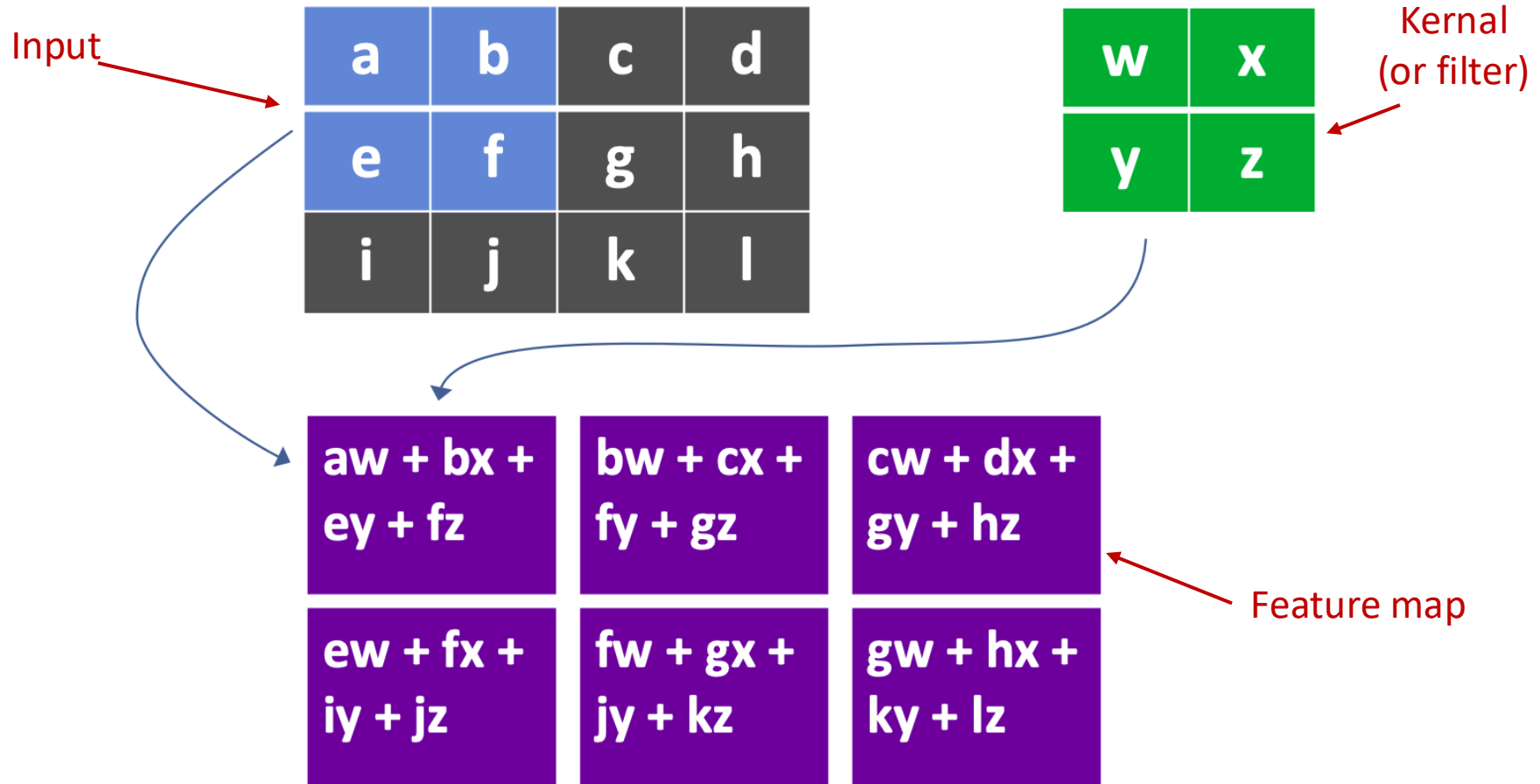


Figure from Balcan (2018)

Advantages of CNNs

- Sparse interaction
 - Reduces memory requirements
 - Improves statistical efficiency
- Parameter sharing
 - The same kernel are used repeatedly
- Equivariant representations
 - transforming the input = transforming the output
 - Useful when care only about the existence of a **pattern**, rather than the **location**

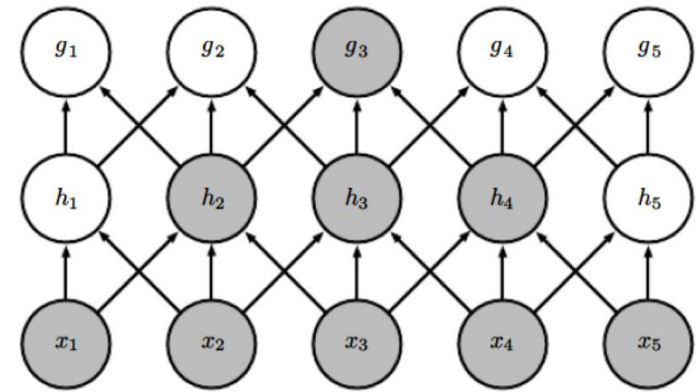
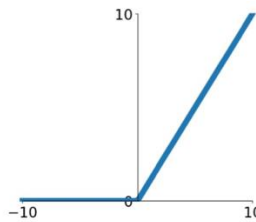


Figure from *Deep Learning*, by Goodfellow, Bengio, and Courville

Other Layers: Activation Functions

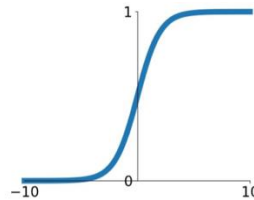
- Activation functions determine whether a neuron is activated based on its input, effectively deciding whether the input is important for making a prediction.

- ReLU: $\sigma(x) = \max(0, x)$



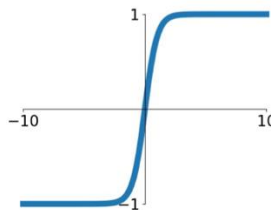
Q: Any other benefit?

- Sigmoid: $\sigma(x) = \frac{1}{1+e^{-x}}$



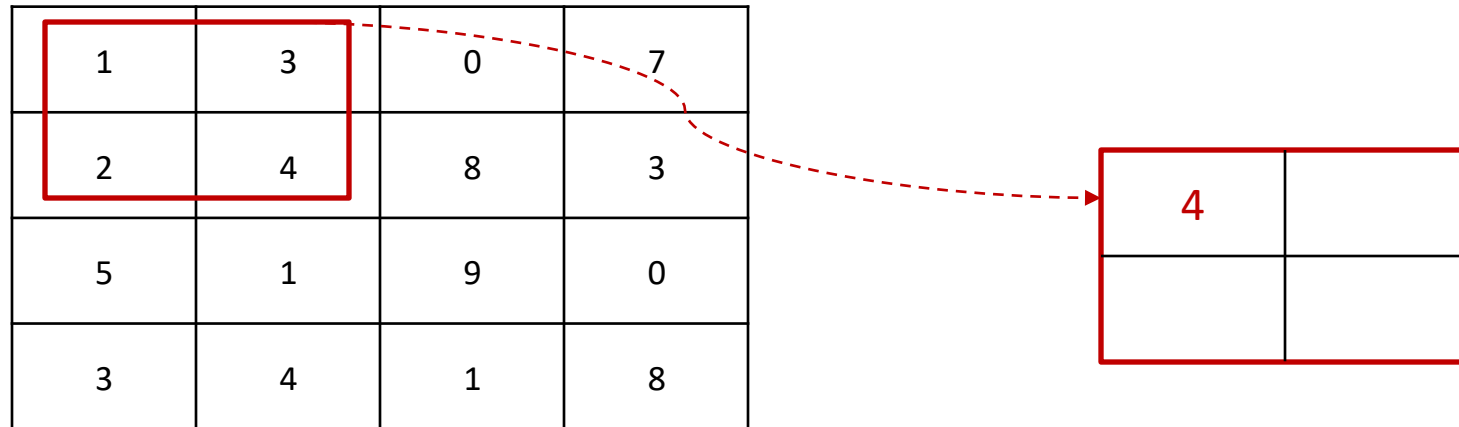
Activation Functions introduce **non-linearity**, enabling the network to learn complex patterns and model intricate relationships within data.

- tanh: $\sigma(x) = \tanh(x)$



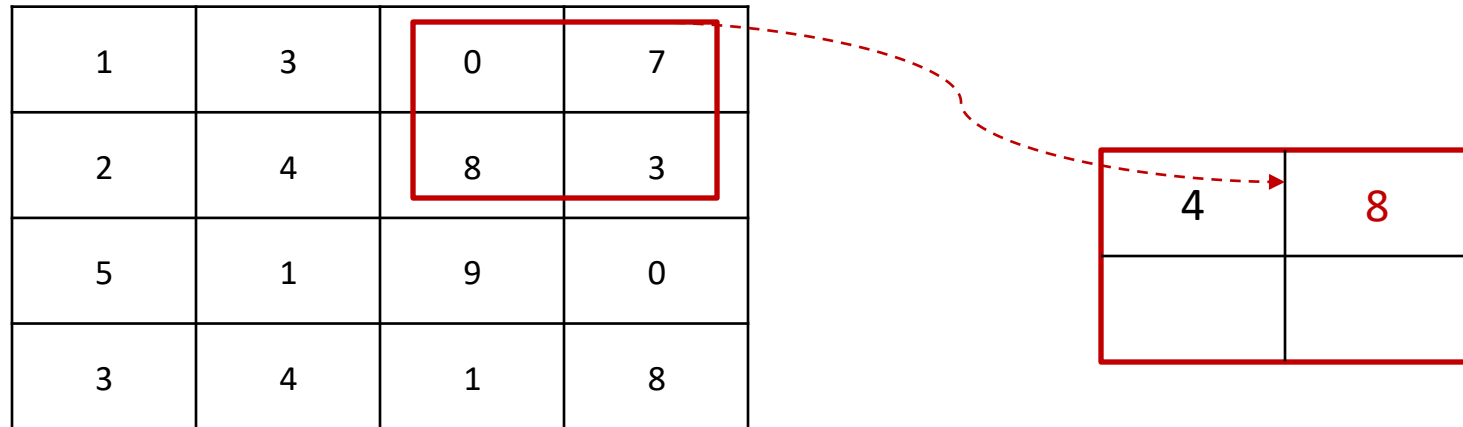
Other Layers: Pooling Layer

- We use a pooling layer to downsize the inputs.
- For example, max pooling (2x2 filter and stride 2)



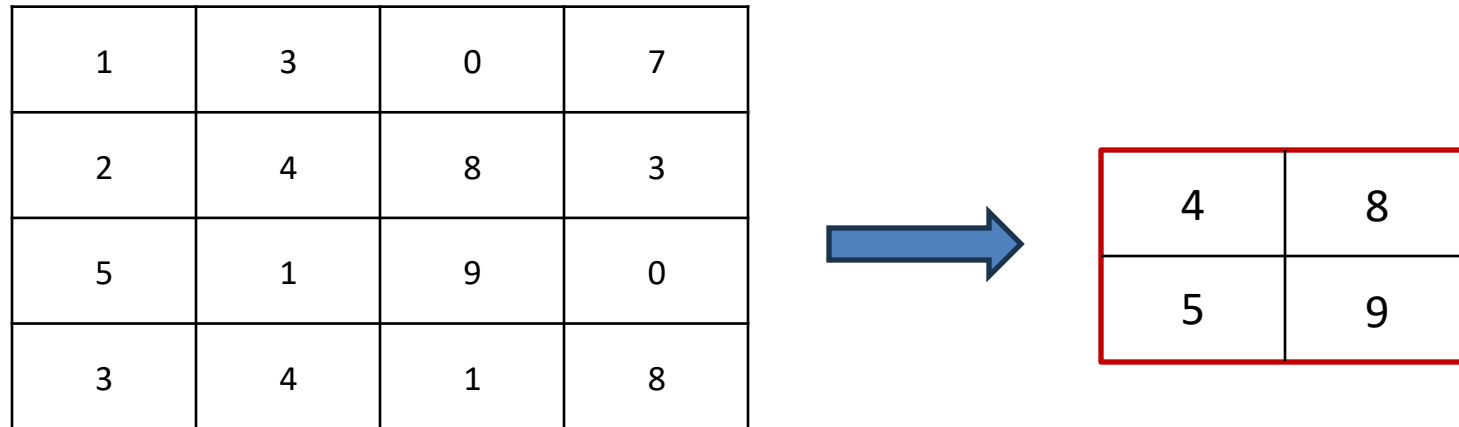
Other Layers: Pooling Layer

- We use a pooling layer to downsize the inputs.
- For example, max pooling (2x2 filter and stride 2)

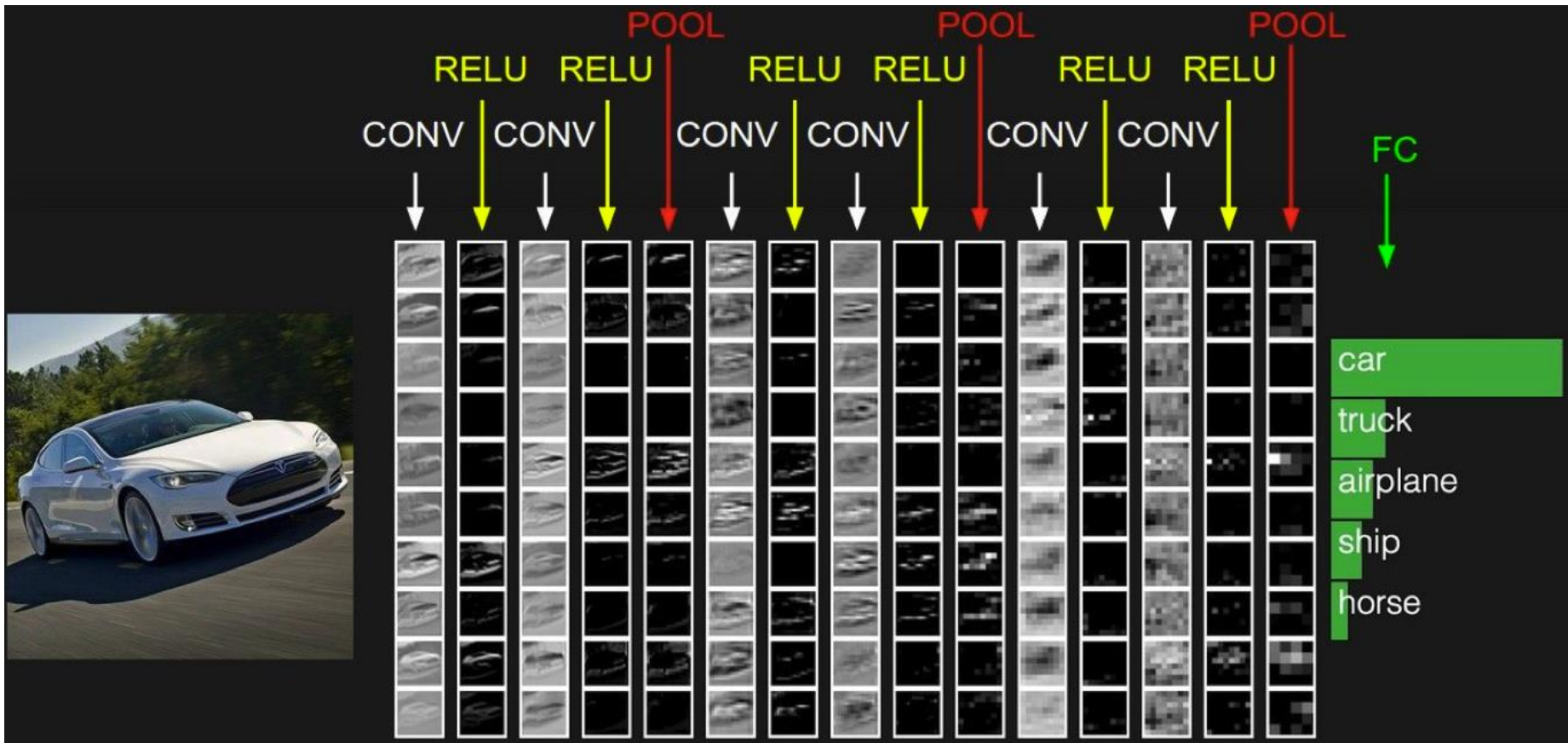


Other Layers: Pooling Layer

- We use a pooling layer to downsize the inputs.
- For example, max pooling (2x2 filter and stride 2)



A Case of CNNs



Q: How to update the parameters of each layer?

Figure from *Feifei Li & Andrej Karpathy (2016)*

