

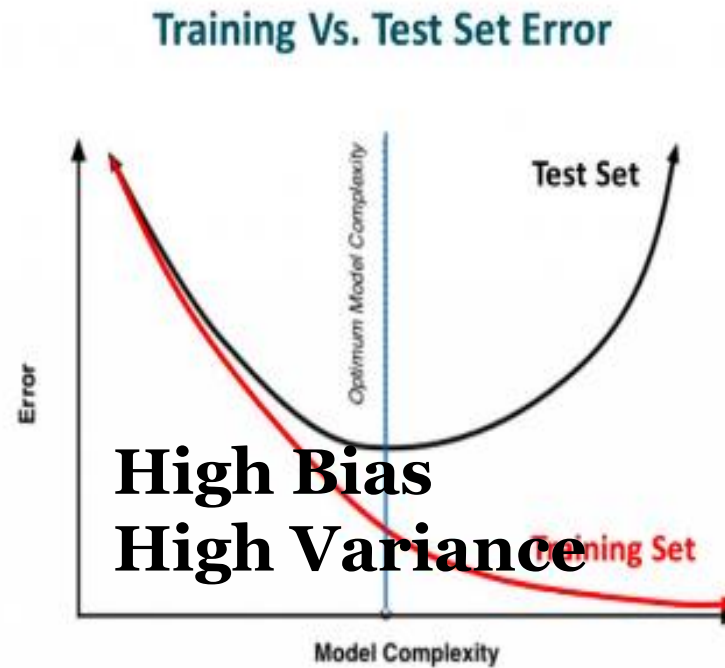
# Deep Learning Regularization and Optimization

Deep Learning  
E. Fatemizadeh Fall 2021



# Regularization

- › From Previous Discussion:
  - The Bias-Variance Tradeoff (Dilemma)



# Regularization

## › The Bias-Variance Tradeoff (Dilemma)

$$y = f(x) + \varepsilon$$

$$\mathbb{E} [(y - \hat{f})^2] = \mathbb{E} [(f + \varepsilon - \hat{f})^2]$$

$$= \mathbb{E} [(f + \varepsilon - \hat{f} + \mathbb{E}[\hat{f}] - \mathbb{E}[\hat{f}])^2]$$

$$= (f - \mathbb{E}[\hat{f}])^2 + \text{Var}[y] + \text{Var} [\hat{f}]$$

$$= \text{Bias}[\hat{f}]^2 + \text{Var}[y] + \text{Var} [\hat{f}]$$

$$= \text{Bias}[\hat{f}]^2 + \sigma^2 + \text{Var} [\hat{f}]$$

Lower Bound on the  
Expected Error on  
**Unseen/Validation**  
Samples



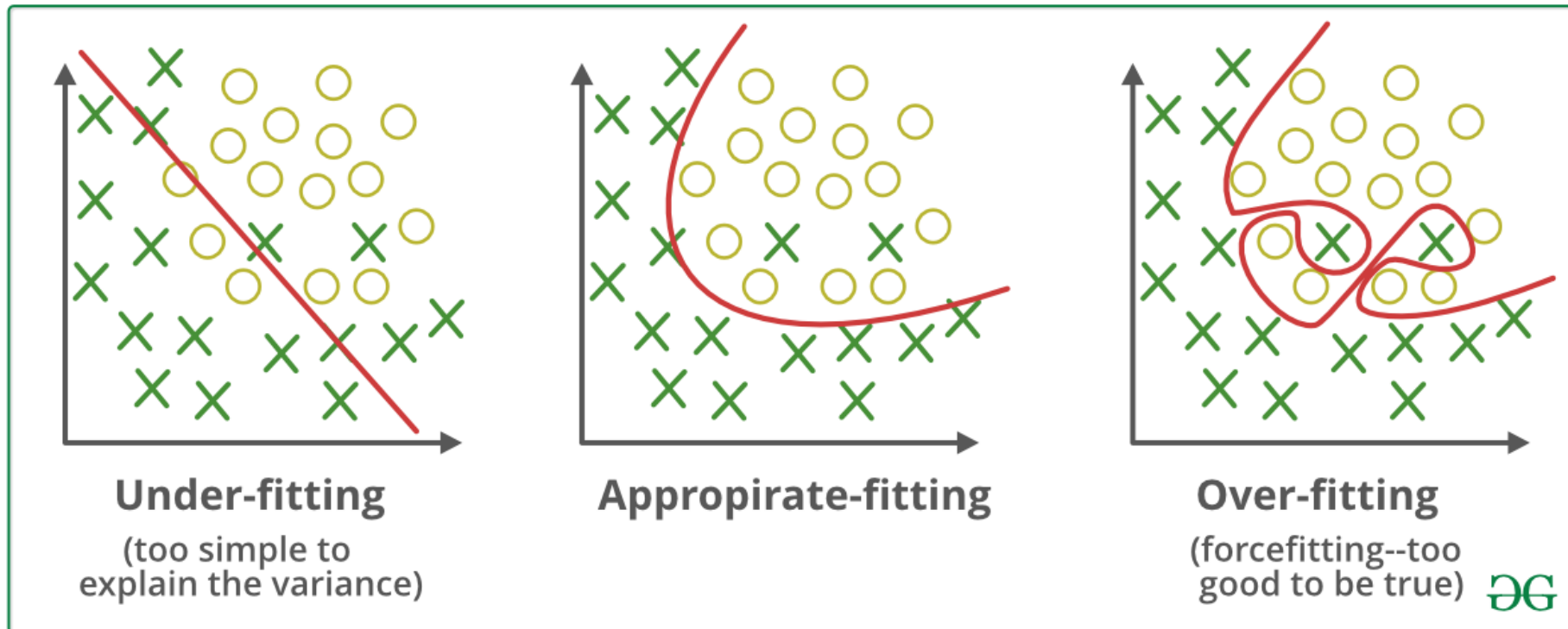
# Regularization

- › The Bias-Variance Tradeoff (Dilemma)
- › Simple models trained on different samples of data **do not differ much** from each other. However, far from the true (underfitting).
- › On the other hand, complex models trained on different subsets of data **are very different** from each other (overfitting).
- › Simple model: high bias, low variance
- › Complex model: low bias, high variance



# Regularization

- › Motivation:
  - Overfitting Challenge



# Regularization

## › Overfitting Challenge (What to do?)

Training Error	Valid Error	Cause	Solution
High	High	High bias	- Increase model complexity - Train for more epochs
Low	High	High variance	- Add more training data (e.g., <b><u>dataset augmentation</u></b> ) - Use <b><u>regularization</u></b> - User <b>early stopping</b> (train less)
Low	Low	Perfect tradeoff	- You are done!



# Regularization

## › Definition:

*“**Any** modification we make to a learning algorithm that is intended to **reduce** its **generalization error** but not its **training error**.”*

## › Solution:

- *Parameter (weights) Penalties*
- *Dataset Augmentation*
- *Noise Robustness (input/output)*
- *Early Stopping*
- *Bagging and Other Ensemble Methods*
- *Dropout*



# Parameter Penalties

› Instead *empirical risk minimization*:

$$\text{minimize}(\text{Loss}(\text{Data}|\text{Model}))$$

› Try *structural risk minimization*:

$$\text{minimize}(\text{Loss}(\text{Data}|\text{Model}) + \lambda \text{Complexity}(\text{Model}))$$

› Two major policies for model complexity::

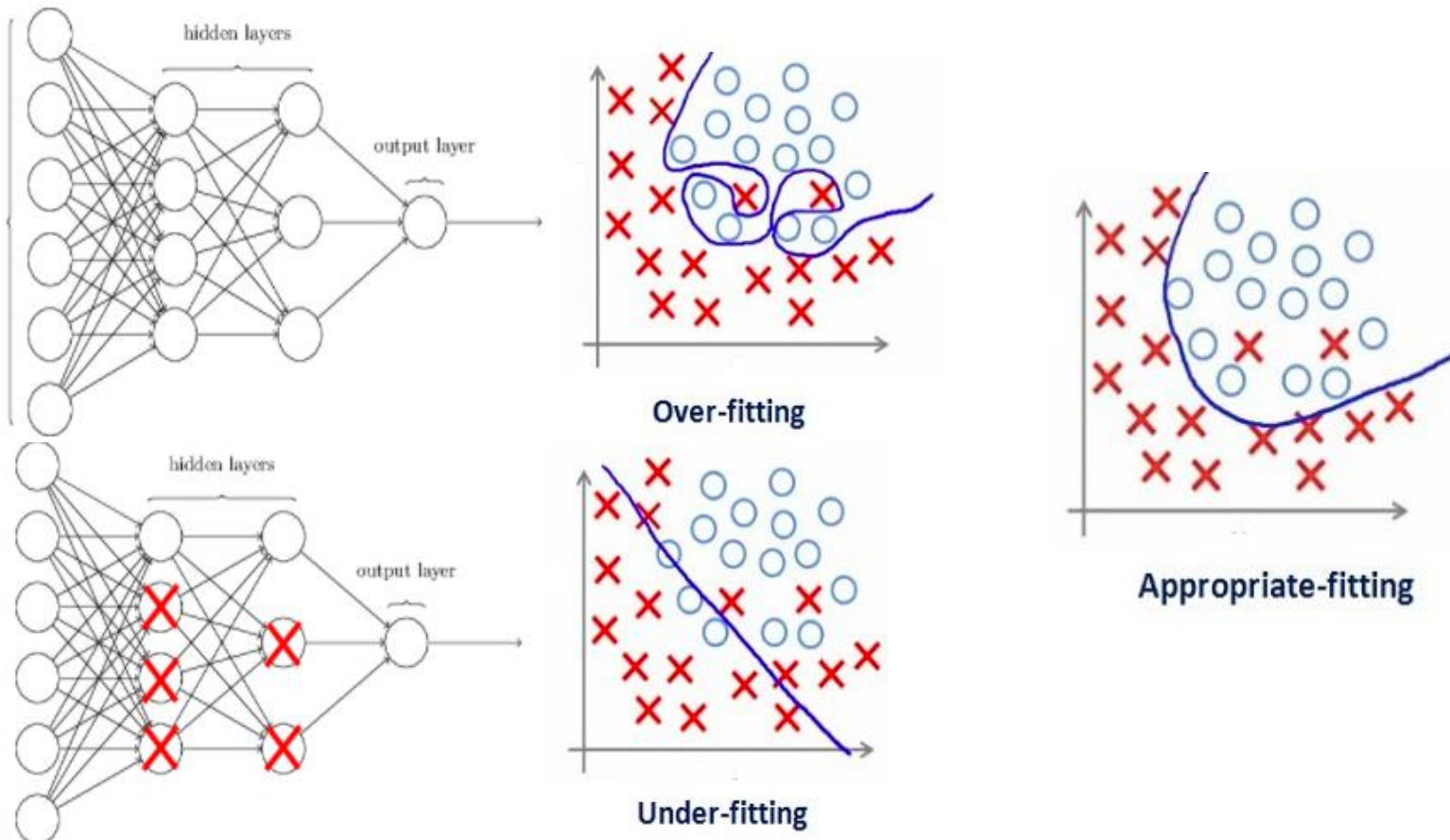
- A function of the weights of all the features in the model.
- A function of the total number of features with nonzero weights.





# Parameter Penalties

› How it works:



# Regularization by Parameters Penalties

› General Formulation:

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha\Omega(\theta)$$

› Typically Not applied on ***Bias*** terms (why)?



# L<sup>2</sup> Regularization

- › An old fashion technique!
- › Idea is **weight decay**
- › Known as:
  - Ridge Regression
  - Tikhonov Regularization

$$J_{L_2} = (MSE / CE) + \lambda \sum_{Weight \setminus Bias} \|w\|_2^2$$

or

$$J_{L_2} = (MSE / CE) + \sum_{Layers} \lambda_{\text{Layers}} \sum_{Weight \setminus Bias} \|w\|_2^2$$



# L<sup>2</sup> Regularization – Deep Learning, Goodfellow 2016

› Formulation:

$$\tilde{J}(w; X, y) = \frac{\alpha}{2} w^\top w + J(w; X, y), \quad (7.2)$$

with the corresponding parameter gradient

$$\nabla_w \tilde{J}(w; X, y) = \alpha w + \nabla_w J(w; X, y). \quad (7.3)$$

To take a single gradient step to update the weights, we perform this update:

$$w \leftarrow w - \epsilon (\alpha w + \nabla_w J(w; X, y)). \quad (7.4)$$

Written another way, the update is:

$$w \leftarrow (1 - \epsilon\alpha)w - \epsilon \nabla_w J(w; X, y). \quad (7.5)$$



# L<sup>2</sup> Regularization - Analysis

› If we know exact solution:

$$w^* = \arg \min_w J(w).$$

› 2<sup>nd</sup> order Approximation:

$$\hat{J}(\theta) = J(w^*) + \frac{1}{2}(w - w^*)^\top H(w - w^*),$$

$$\nabla_w \hat{J}(w) = H(w - w^*)$$

› With L<sup>2</sup> term:

$$\alpha \tilde{w} + H(\tilde{w} - w^*) = 0$$

$$(H + \alpha I)\tilde{w} = Hw^*$$

$$\tilde{w} = (H + \alpha I)^{-1} Hw^*.$$



# L<sup>2</sup> Regularization - Analysis

- › For Convex Loss function,  $H$ , is pdm (Unitary  $Q$  and positive  $\Lambda$ ) :

$$H = Q\Lambda Q^\top.$$

$$\begin{aligned}\tilde{w} &= (Q\Lambda Q^\top + \alpha I)^{-1} Q\Lambda Q^\top w^* \\ &= \left[ Q(\Lambda + \alpha I)Q^\top \right]^{-1} Q\Lambda Q^\top w^* \\ &= Q(\Lambda + \alpha I)^{-1} \Lambda Q^\top w^*.\end{aligned}$$



# L<sup>2</sup> Regularization - Analysis

› For Convex Loss function, H, is pdm, and :

$$(\Lambda + \alpha I)^{-1} \Lambda \iff \frac{\lambda_i}{\lambda_i + \alpha}$$

› Numerical Examples:

$$H = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix} \Rightarrow Q \Lambda^{-1} \Lambda Q^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$Q(\Lambda + 2I)^{-1} \Lambda Q^T = \begin{bmatrix} 0.46 & -0.14 & -0.04 \\ -0.14 & 0.43 & -0.14 \\ -0.04 & -0.14 & 0.46 \end{bmatrix}$$

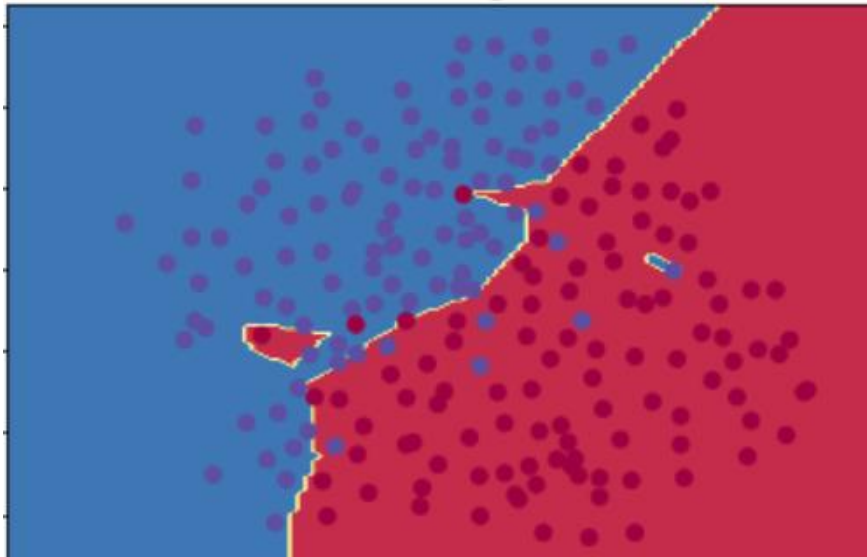


# L<sup>2</sup> Regularization - Analysis

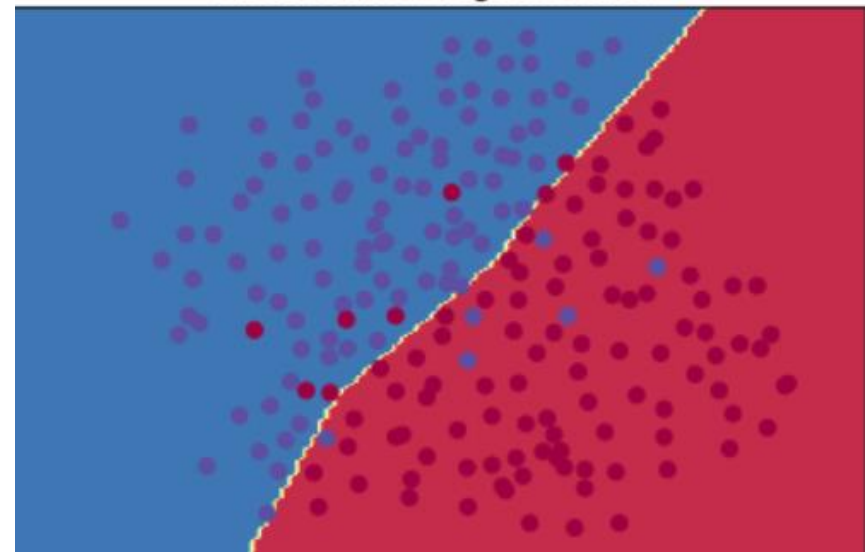
› For Convex Loss function, H, is pdm, and :

$$(\Lambda + \alpha I)^{-1} \Lambda \iff \frac{\lambda_i}{\lambda_i + \alpha}$$

Model without regularization



Model with L2-regularization





# $L^2$ Regularization - Analysis

› In Linear Regression:

$$(Xw - y)^\top (Xw - y). \quad (7.14)$$

When we add  $L^2$  regularization, the objective function changes to

$$(Xw - y)^\top (Xw - y) + \frac{1}{2}\alpha w^\top w. \quad (7.15)$$

This changes the normal equations for the solution from

$$w = (X^\top X)^{-1} X^\top y \quad (7.16)$$

to

$$w = (X^\top X + \alpha I)^{-1} X^\top y. \quad (7.17)$$



# L<sup>1</sup> Regularization

› Idea is **weight Sparsification**:

$$\tilde{J}(w; X, y) = \alpha ||w||_1 + J(w; X, y),$$

$$\nabla_w \tilde{J}(w; X, y) = \alpha \text{sign}(w) + \nabla_w J(X, y; w)$$

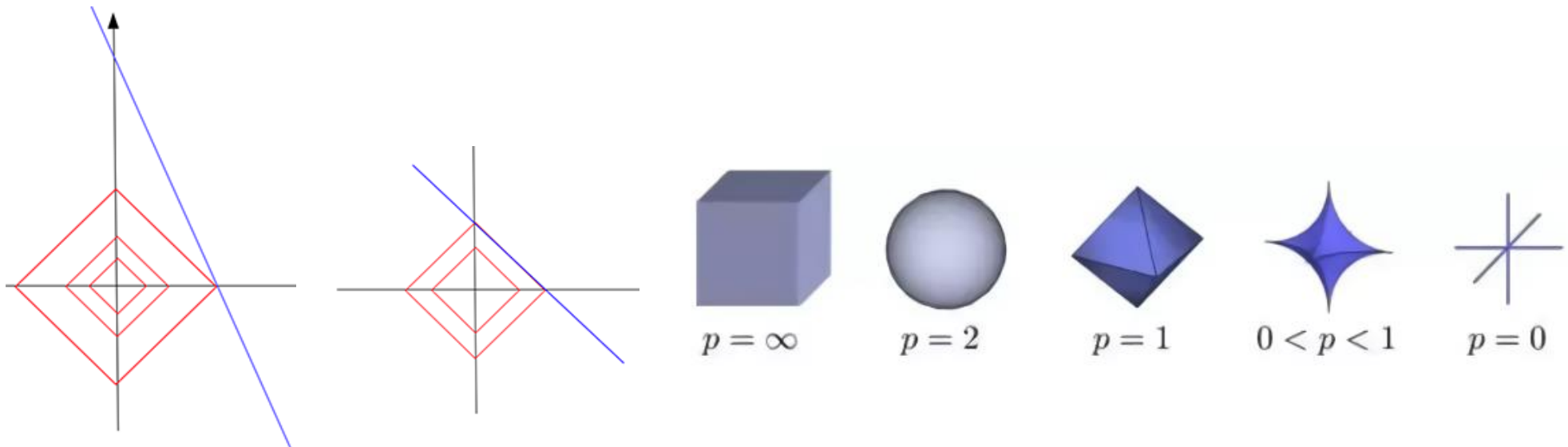
$$\nabla_w \hat{J}(w) = H(w - w^*),$$



# $L^1$ Regularization

## › Why Sparsify?

- Suppose we want solve  $AX=b$  with minimizing  $L_1$ 
  - › May have infinite solution!!!



# L<sup>1</sup> Regularization

› For diagonal Hessian Matrix:

$$\hat{J}(w; X, y) = J(w^*; X, y) + \sum_i \left[ \frac{1}{2} H_{i,i} (w_i - w_i^*)^2 + \alpha |w_i| \right].$$

$$w_i = \text{sign}(w_i^*) \max \left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}.$$



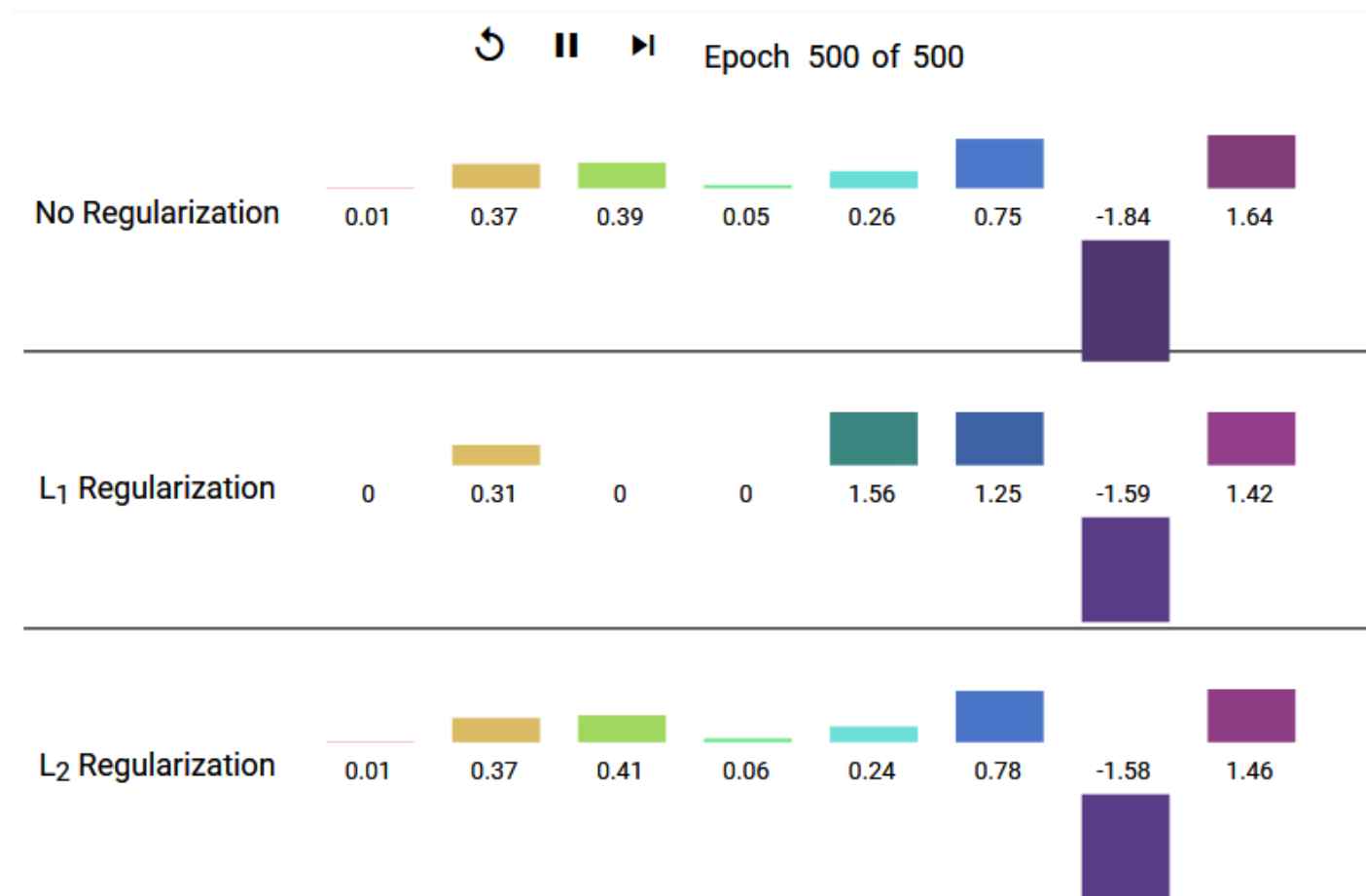
# Regularization Effect

<https://developers.google.com/machine-learning/crash-course/>



# Regularization Effect

<https://developers.google.com/machine-learning/crash-course/>



# Elastic Net Regularization - Regularization and Variable Selection via the Elastic Net, Zou and Hastie, 2004

›  $L_1 + L_2$  Regularization:

$$J_{Elastic} = (MSE / CE) + \lambda \left( \alpha \sum_{Weight \setminus Bias} \|w\|_2^2 + (1 - \alpha) \sum_{Weight \setminus Bias} \|w\|_1 \right)$$



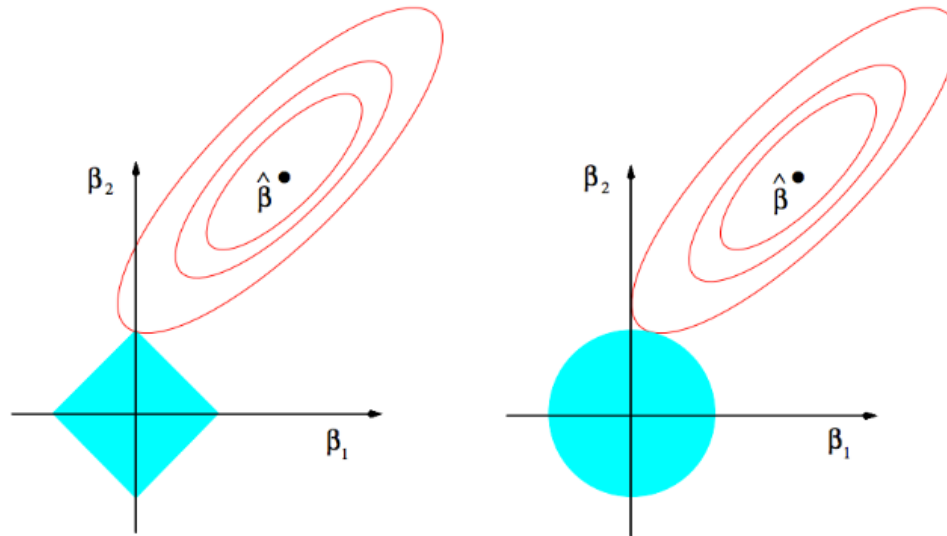
# Geometric Perspective

› The Problems:

$$\text{minimize} \left( \text{Loss}(\text{Data}|\text{Model}) + \lambda \text{Complexity}(\text{Model}) \right)$$

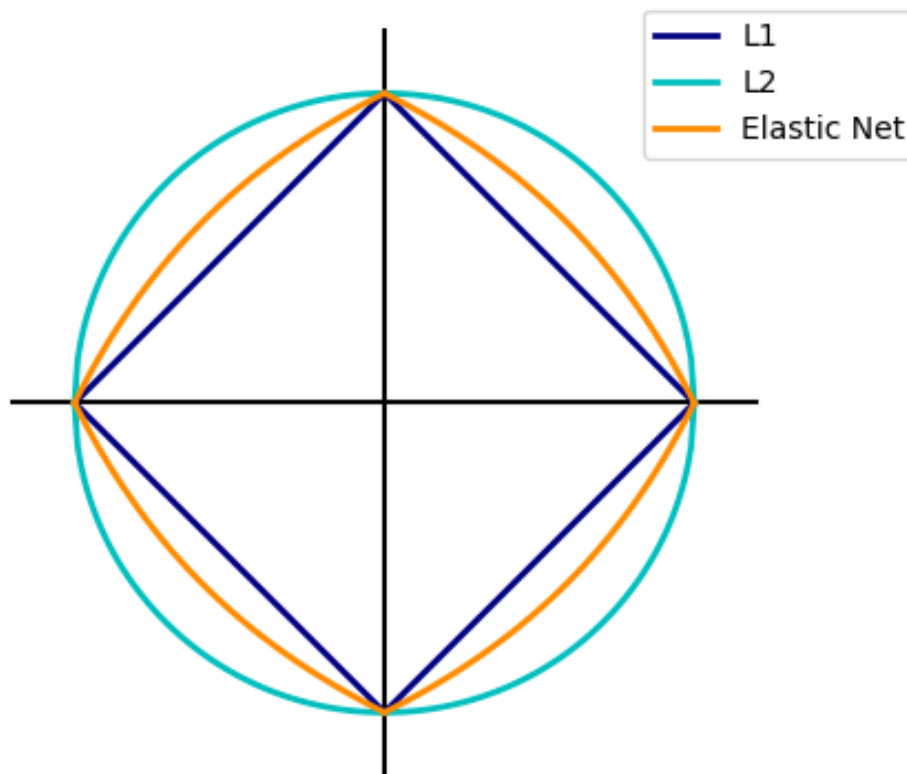
› **Is equivalent to:**

$$\text{minimize} \left( \text{Loss}(\text{Data}|\text{Model}), \quad \text{Complexity}(\text{Model}) \right) \leq \eta$$





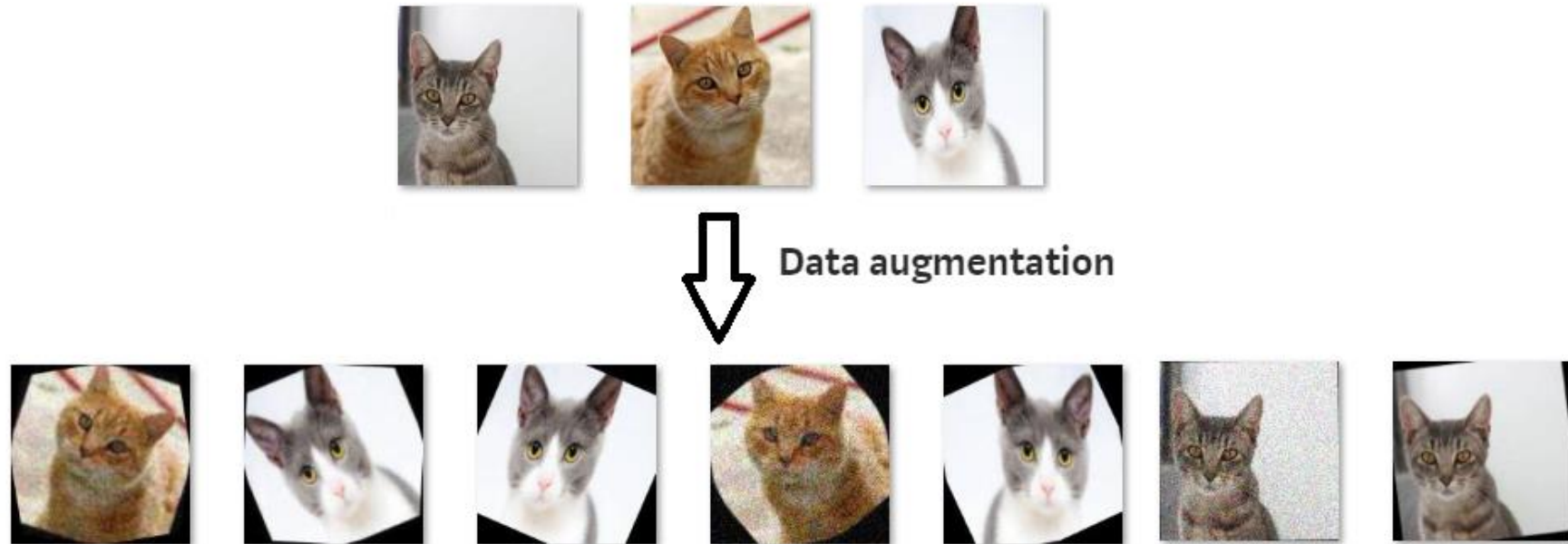
# Geometric Perspective



# Data Augmentation


## › Idea:

- Invariance (Rotation, Scaling, Translation, Shearing, Mirroring, ...) and Regularization



# Data Augmentation

## › Standard Transform

- Horizontal/Vertical Flip (Mirroring)
- Horizontal/Vertical Translation
- Horizontal/Vertical Scale
- Rotation
- Shearing (Skewing) 
- Contrast/Brightness/illumination/Color/... Change
- Random Crop and then Scale
- Random Cutout (Blackout a random rectangular/square)
- Additive/Multiplicative Noise

## › Modern Approach:

- Deep Generated Data, Augmented Network (DeepFake, ...)



# Regularization – Dropout - Dropout: A Simple Way to Prevent Neural Networks from Overfitting, Srivastava, 2014

## › Preliminary:

- Ensemble Methods: Learn multiple models (Weak Classifier) to solve the same problem and then combine (aggregate) them for better results
- Combination of “Different/Same”  
“Model/Dataset/Features/Objective/ Optimization/Initials”

## › Policies:

- Bagging
- Boosting
- Stacking



# Regularization – Dropout

- › Bagging (Bootstrap aggregating):
  1.  $K$  new training datasets are generated by uniformly random sampling with **replacement (Bootstrapping)** from the original dataset.
  2. Train  $K$  model independently
  3. Combine/Aggregate the classifier/regressor using:
    1. Averaging (Arithmetic, Geometric, ...)
    2. Voting
- › Boosting: Each data and classifier has its own weight (determined with boosting procedure), AdaBoost is most known methods.
- › Stacking: Train  $K$  models using all data, then train a combiner/aggregator classifier to combine  $K$  result, as best as possible.



# Regularization – Dropout

## › Motivation:

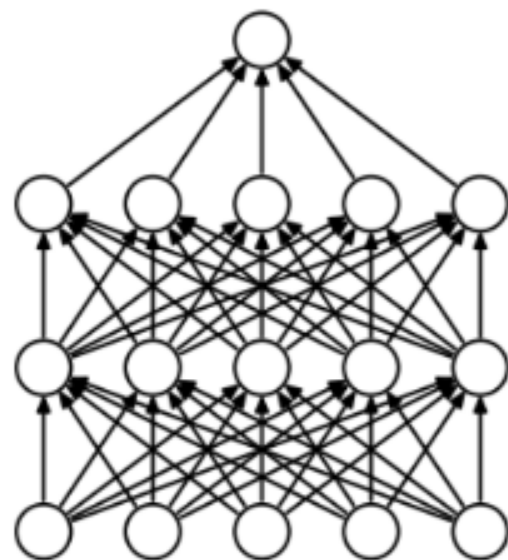
### – Co-Adaptation:

*If two or more neurons extract the same feature repeatedly or highly correlated behavior (co-adaptation), the network isn't reaching its full capability.*

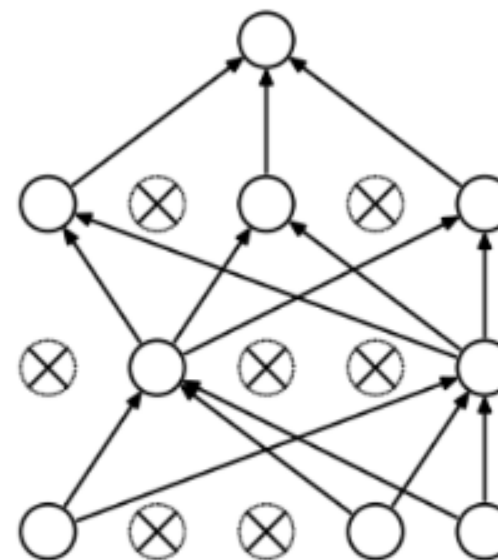


# Regularization – Dropout

- › Solution: Update (EBP) a fraction (Randomly Chosen) of all weights is each iteration!



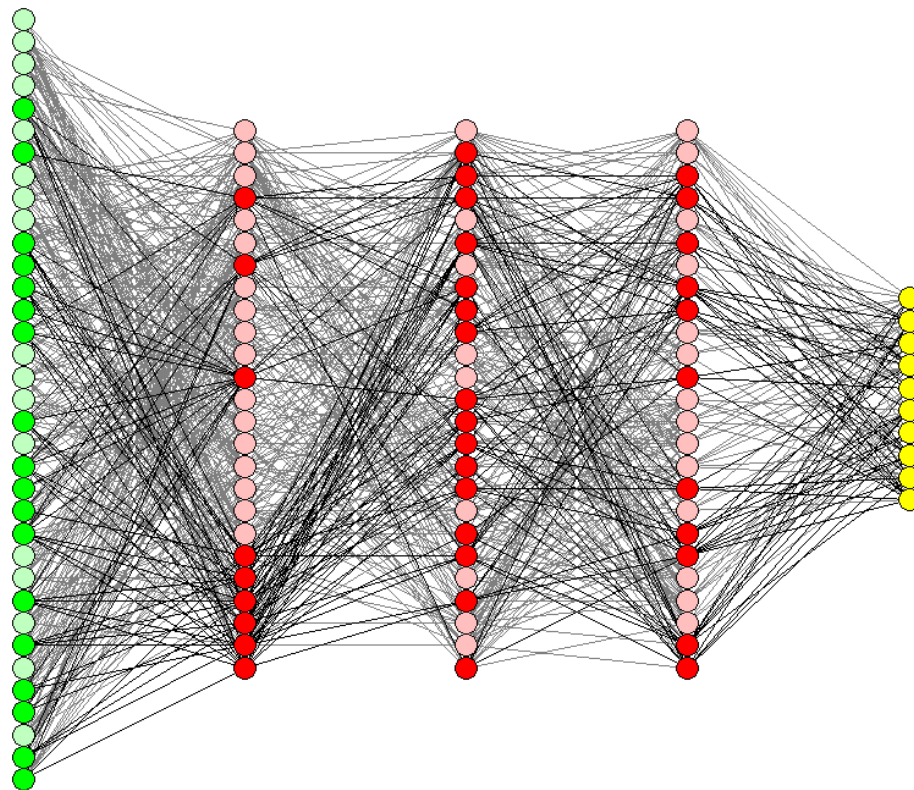
(a) Standard Neural Net



(b) After applying dropout.



# Regularization – Dropout





# Regularization – Dropout

## › Training Phase:

- For each **hidden layer**, for each **training sample**, for each **iteration**, dropout (zero out) a random fraction,  $(1-p)$ , of neuron (and corresponding weight).
  - › Input Layer:  $P > 0.8$  or  $P = 1.0$
  - › Hidden Layer:  $P = 0.5$
  - › Output Layer:  $P = 1.0$

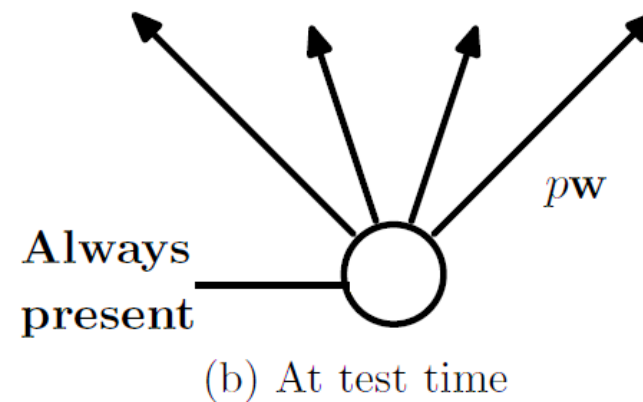
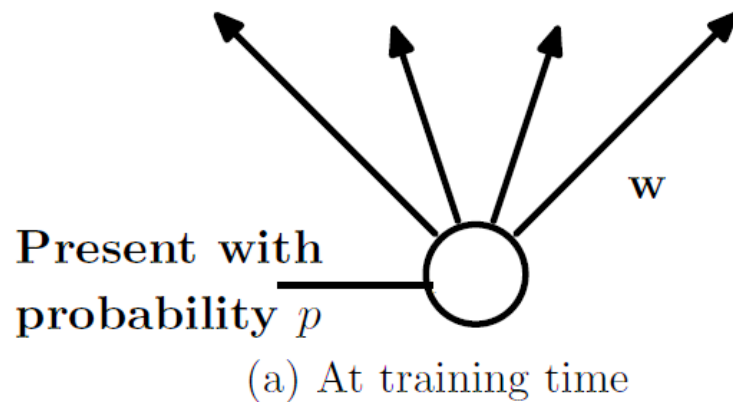
## › Test Phase:

- Use all activations, but reduce them by a factor  $(1-p)$  (to account for the missing activations during training).



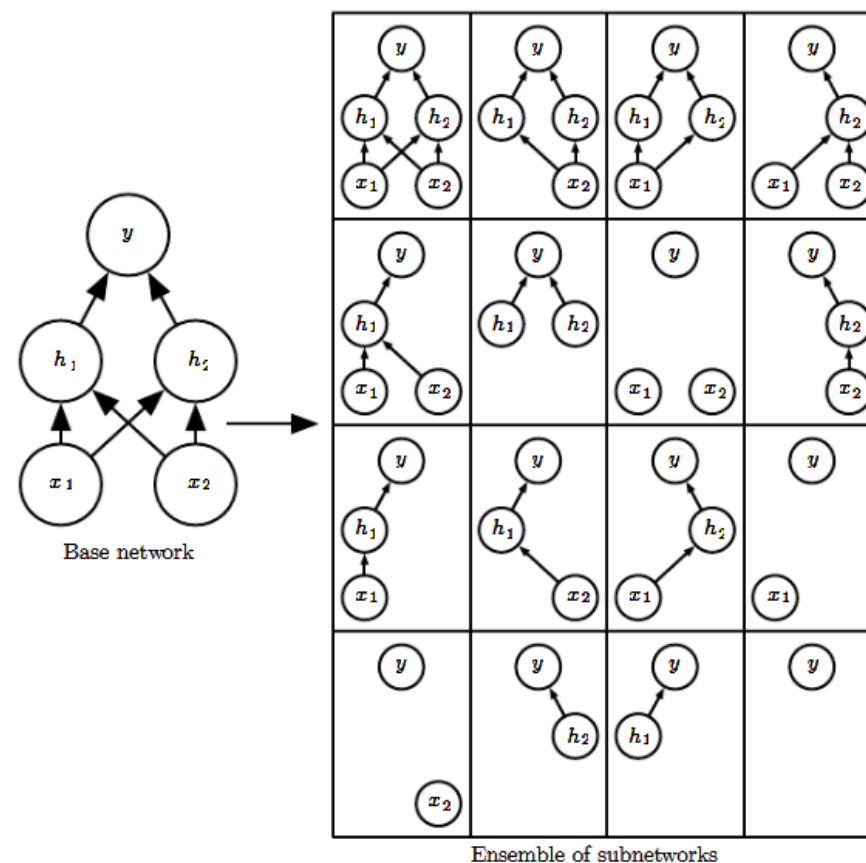
# Regularization – Dropout

## › Dropout Training/Test Scheme



# Regularization – Dropout

## › Dropout Training Scheme



# Regularization – Dropout. [Understanding Dropout, Baldi, NIPS2013](#)

- › Mathematics of Dropout (N: Complete NN, D: Dropout NN)
- › Dropping out is modelled by Bernoulli distribution:

$$\delta_i \sim \text{Bernoulli}(p)$$

$$E_N = \frac{1}{2} \left( t - \sum_{i=1}^n w'_i I_i \right)^2$$

$$E_D = \frac{1}{2} \left( t - \sum_{i=1}^n \delta_i w_i I_i \right)^2$$

$$E_N = \frac{1}{2} \left( t - \sum_{i=1}^n p_i w_i I_i \right)^2$$



# Regularization – Dropout. [Understanding Dropout, Baldi, NIPS2013](#)

- › Mathematics of Dropout (N: Complete NN, D: Dropout NN)

$$\frac{\partial E_D}{\partial w_i} = -t\delta_i I_i + w_i \delta_i^2 I_i^2 + \sum_{j=1, j \neq i}^n w_j \delta_i \delta_j I_i I_j$$

$$\frac{\partial E_N}{\partial w_i} = -tp_i I_i + w_i p_i^2 I_i^2 + \sum_{j=1, j \neq i}^n w_j p_i p_j I_i I_j$$

$$E \left[ \frac{\partial E_D}{\partial w_i} \right] = -tp_i I_i + w_i p_i^2 I_i^2 + w_i \text{Var}(\delta_i) I_i^2 + \sum_{j=1, j \neq i}^n w_j p_i p_j I_i I_j$$

$$= \frac{\partial E_N}{\partial w_i} + w_i \text{Var}(\delta_i) I_i^2$$

$$= \frac{\partial E_N}{\partial w_i} + w_i p_i (1 - p_i) I_i^2$$



# Regularization – Dropout. [Understanding Dropout, Baldi, NIPS2013](#)

› It is Like Regularization

$$\begin{aligned} E\left[\frac{\partial E_D}{\partial w_i}\right] &= -tp_i I_i + w_i p_i^2 I_i^2 + w_i \text{Var}(\delta_i) I_i^2 + \sum_{j=1, j \neq i}^n w_j p_i p_j I_i I_j \\ &= \frac{\partial E_N}{\partial w_i} + w_i \text{Var}(\delta_i) I_i^2 \\ &= \frac{\partial E_N}{\partial w_i} + \boxed{w_i p_i (1 - p_i) I_i^2} \end{aligned}$$

$$E_R = \frac{1}{2} \left( t - \sum_{i=1}^n p_i w_i I_i \right)^2 + \sum_{i=1}^n p_i (1 - p_i) w_i^2 I_i^2$$

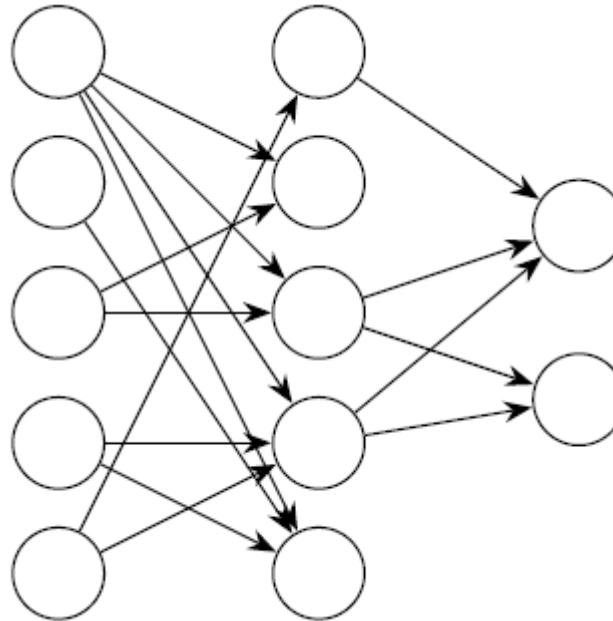
› Why  $p=0.5$  is recommended?



# Regularization – Dropout

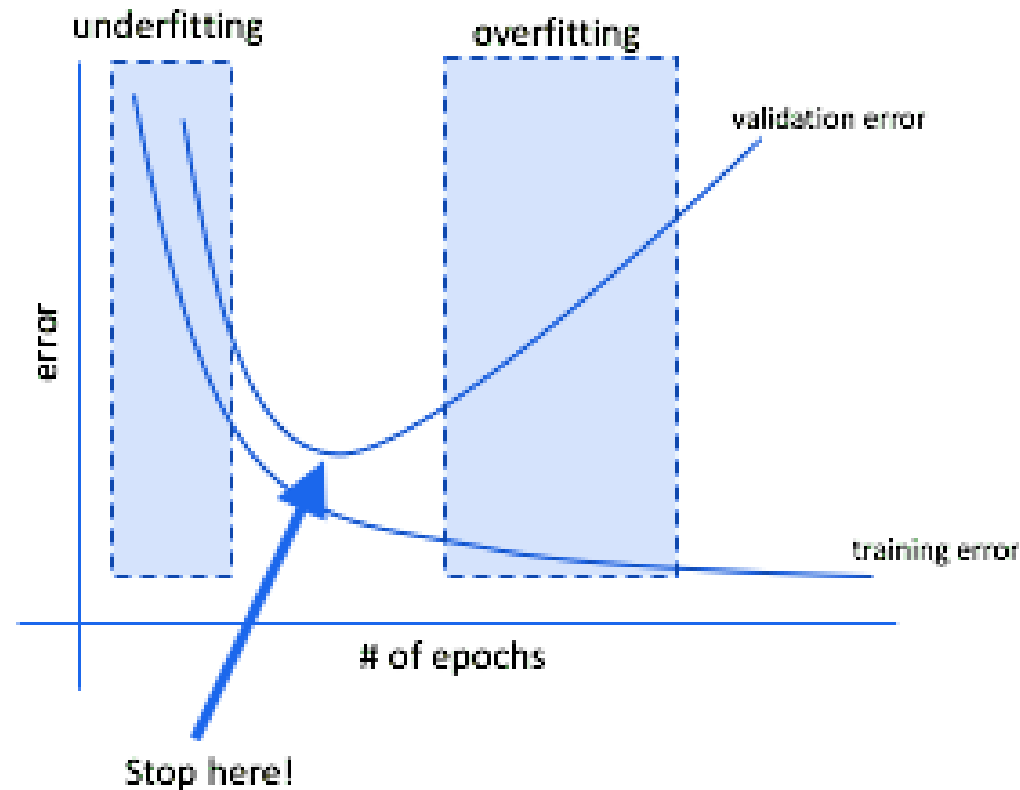
## › Some Variations:

- Gaussian Dropout Multiplicative:  $\delta_i \sim N(1, \sigma^2)$
- Gaussian Dropout Additive:  $\delta_i \sim N(0, \sigma^2)$
- DropConnect:



# Regularization - Early Stopping

› Again Validation Error Curve!





# Regularization - Early Stopping Meta Algorithm

Let  $n$  be the number of steps between evaluations.

Let  $p$  be the “patience,” the number of times to observe worsening validation set error before giving up.

Let  $\theta_o$  be the initial parameters.

$\theta \leftarrow \theta_o$

$i \leftarrow 0$

$j \leftarrow 0$

$v \leftarrow \infty$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

**while**  $j < p$  **do**

    Update  $\theta$  by running the training algorithm for  $n$  steps.

$i \leftarrow i + n$

$v' \leftarrow \text{ValidationSetError}(\theta)$

**if**  $v' < v$  **then**

$j \leftarrow 0$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

$v \leftarrow v'$

**else**

$j \leftarrow j + 1$

**end if**

**end while**

Best parameters are  $\theta^*$ , best number of training steps is  $i^*$ .

---



# Regularization - Early Stopping How to use data #1!

---

**Algorithm 7.2** A meta-algorithm for using early stopping to determine how long to train, then retraining on all the data.

---

Let  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  be the training set.

Split  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  into  $(\mathbf{X}^{(\text{subtrain})}, \mathbf{X}^{(\text{valid})})$  and  $(\mathbf{y}^{(\text{subtrain})}, \mathbf{y}^{(\text{valid})})$  respectively.

Run early stopping (algorithm 7.1) starting from random  $\theta$  using  $\mathbf{X}^{(\text{subtrain})}$  and  $\mathbf{y}^{(\text{subtrain})}$  for training data and  $\mathbf{X}^{(\text{valid})}$  and  $\mathbf{y}^{(\text{valid})}$  for validation data. This returns  $i^*$ , the optimal number of steps.

Set  $\theta$  to random values again.

Train on  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  for  $i^*$  steps.

---



# Regularization - Early Stopping How to use data #2!

---

**Algorithm 7.3** Meta-algorithm using early stopping to determine at what objective value we start to overfit, then continue training until that value is reached.

---

Let  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  be the training set.

Split  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  into  $(\mathbf{X}^{(\text{subtrain})}, \mathbf{X}^{(\text{valid})})$  and  $(\mathbf{y}^{(\text{subtrain})}, \mathbf{y}^{(\text{valid})})$  respectively.

Run early stopping (algorithm 7.1) starting from random  $\theta$  using  $\mathbf{X}^{(\text{subtrain})}$  and  $\mathbf{y}^{(\text{subtrain})}$  for training data and  $\mathbf{X}^{(\text{valid})}$  and  $\mathbf{y}^{(\text{valid})}$  for validation data. This updates  $\theta$ .

$\epsilon \leftarrow J(\theta, \mathbf{X}^{(\text{subtrain})}, \mathbf{y}^{(\text{subtrain})})$

**while**  $J(\theta, \mathbf{X}^{(\text{valid})}, \mathbf{y}^{(\text{valid})}) > \epsilon$  **do**

    Train on  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  for  $n$  steps.

**end while**

---



# Regularization - Early Stopping Theory

- › Why it Early Stopping acts as regulator?
- › Taylor approximation around true solution ( $w^*$ ), and convex assumption!

$$\hat{J}(\theta) = J(w^*) + \frac{1}{2}(w - w^*)^\top H(w - w^*),$$

$$\nabla_w \hat{J}(w) = H(w - w^*).$$

- › Gradient Descent, starting from ( $w^{(0)}=0$ ,  $\tau$  is iteration counter)

$$\begin{aligned} w^{(\tau)} &= w^{(\tau-1)} - \epsilon \nabla_w \hat{J}(w^{(\tau-1)}) \\ &= w^{(\tau-1)} - \epsilon H(w^{(\tau-1)} - w^*), \end{aligned}$$

$$w^{(\tau)} - w^* = (I - \epsilon H)(w^{(\tau-1)} - w^*).$$



# Regularization - Early Stopping Theory

- › Eigen decomposition of  $\mathbf{H}$ :

$$\begin{aligned} \mathbf{w}^{(\tau)} - \mathbf{w}^* &= (\mathbf{I} - \epsilon \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top)(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*) \\ \mathbf{Q}^\top(\mathbf{w}^{(\tau)} - \mathbf{w}^*) &= (\mathbf{I} - \epsilon \mathbf{\Lambda}) \mathbf{Q}^\top(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*) \end{aligned}$$

- › 1<sup>st</sup> order DE, starting  $\mathbf{w}^{(0)} = \mathbf{0}$ , and small step size ( $|1 - \epsilon \lambda_i| < 1$ )

$$\mathbf{Q}^\top \mathbf{w}^{(\tau)} = [\mathbf{I} - (\mathbf{I} - \epsilon \mathbf{\Lambda})^\tau] \mathbf{Q}^\top \mathbf{w}^*.$$

- › Remember similar analysis in  $L_2$  Regularization ([Slide #14](#)):

$$\frac{\lambda_i}{\lambda_i + \alpha} = 1 - \frac{\alpha}{\lambda_i + \alpha} \quad \longrightarrow \quad \begin{aligned} \mathbf{Q}^\top \tilde{\mathbf{w}} &= (\mathbf{\Lambda} + \alpha \mathbf{I})^{-1} \mathbf{\Lambda} \mathbf{Q}^\top \mathbf{w}^*, \\ \mathbf{Q}^\top \tilde{\mathbf{w}} &= [\mathbf{I} - (\mathbf{\Lambda} + \alpha \mathbf{I})^{-1} \alpha] \mathbf{Q}^\top \mathbf{w}^*. \end{aligned}$$



# Regularization - Early Stopping Theory

› Compare Two approaches:

$$Q^T \tilde{w} = [I - (\Lambda + \alpha I)^{-1} \alpha] Q^T w^*.$$

› Gives:  $Q^T w^{(\tau)} = [I - (I - \epsilon \Lambda)^\tau] Q^T w^*.$

$$(I - \epsilon \Lambda)^\tau = (\Lambda + \alpha I)^{-1} \alpha,$$

› Solve for small  $\epsilon$  and Large  $\alpha$  (Regularization factor):

$$I - \epsilon \tau \Lambda \approx (-\Lambda + \alpha^{-1} I) \alpha = I - \alpha^{-1} \Lambda \rightarrow \alpha \approx \epsilon \tau$$



## Regularization - Others

- › Noise Robustness (adding noise to input/output)
- › Multitasking Learning:
  - Learn two or more tasks with a layer of shared parameters
- › Parameter Tying:
  - Two Model performs same classification task but different input pdfs, we consider a penalty as:

$$\hat{y}^{(A)} = f(\mathbf{w}^{(A)}, \mathbf{x})$$

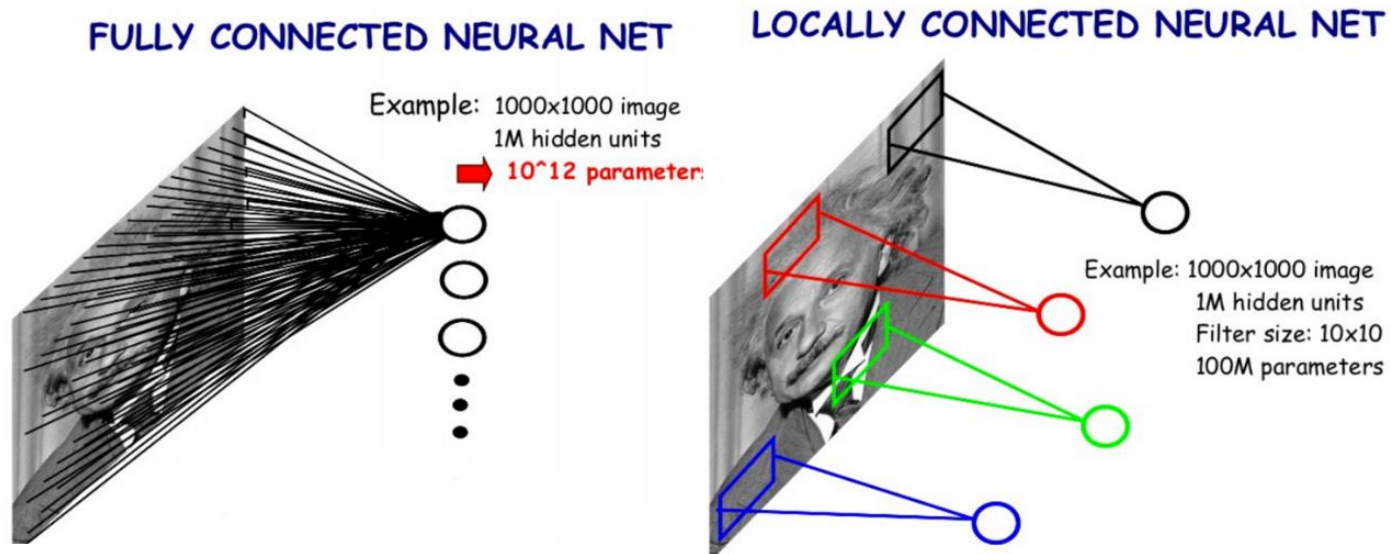
$$\hat{y}^{(B)} = g(\mathbf{w}^{(B)}, \mathbf{x})$$

$$\Omega(\mathbf{w}^{(A)}, \mathbf{w}^{(B)}) = \|\mathbf{w}^{(A)} - \mathbf{w}^{(B)}\|_2^2$$



# Regularization - Others

- › Parameter Sharing:
  - Force two set of weights is equals
- › Sparse Parameters:
  - Force some weights equal to zero





# Optimization Challenges

- › Local minimum and saddle points:
  - For many **high-dimensional non-convex functions** ( and random function) Local minimum/maximum are **rare** compared to saddle points. (#Saddles/#Locals grows exponentially with  $R^D$ )
  - In local minimum, eigenvalues of Hessian are Positive
  - In Saddle Points, eigenvalues of Hessian may be Positive or Negative

- › Ill-Conditioning:

- Hessian term may exceeds gradient term

$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}.$$



# Optimization Challenges

- › Vanishing/Exploding Gradient:
  - Suppose a simple Computational Graph (Back Propagation) in which we repeatedly **multiply** by a matrix  $\mathbf{W}$

$$\mathbf{W} = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}.$$

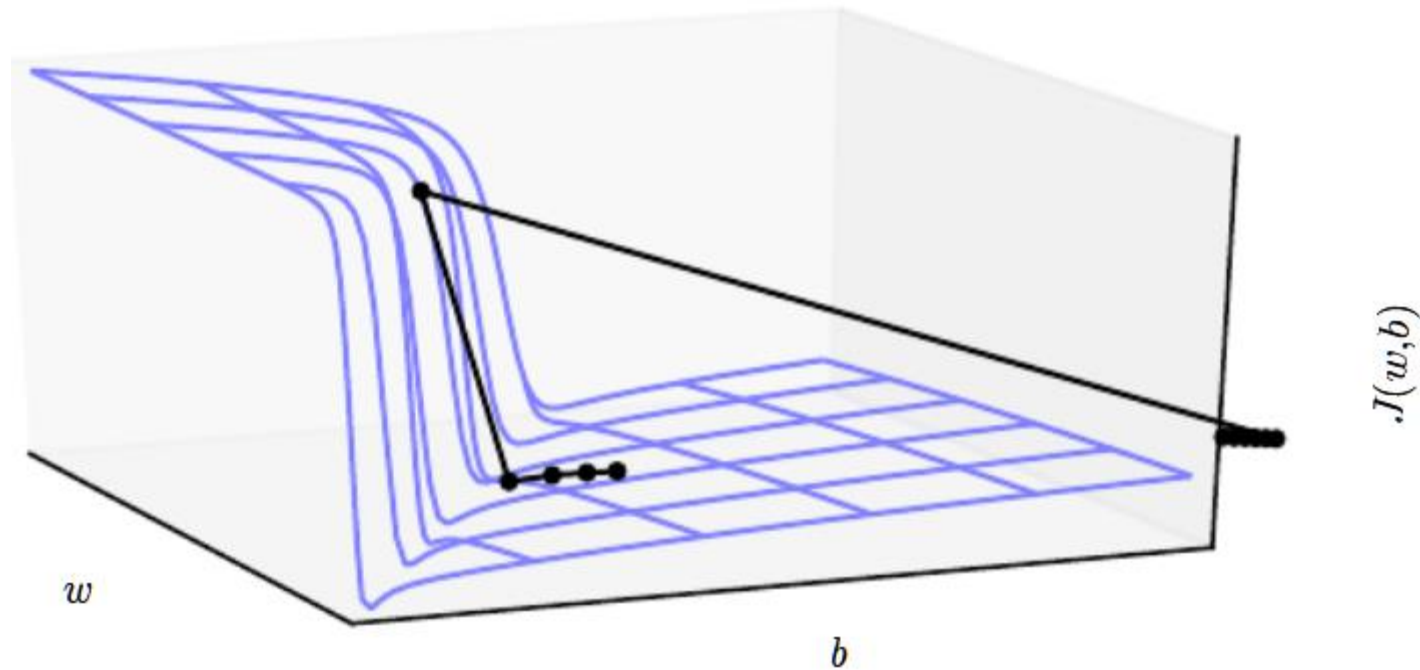
$$\mathbf{W}^t = (\mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1})^t = \mathbf{V} \text{diag}(\boldsymbol{\lambda})^t \mathbf{V}^{-1}.$$

- › Vanishing ( $\lambda_i < 1$ ) or Exploding ( $\lambda_i > 1$ )
  - Vanishing: which direction is?
  - Exploding: Cliff Structure



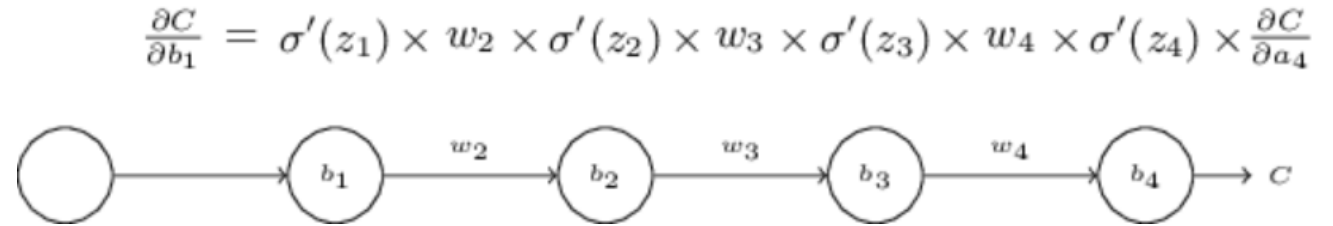
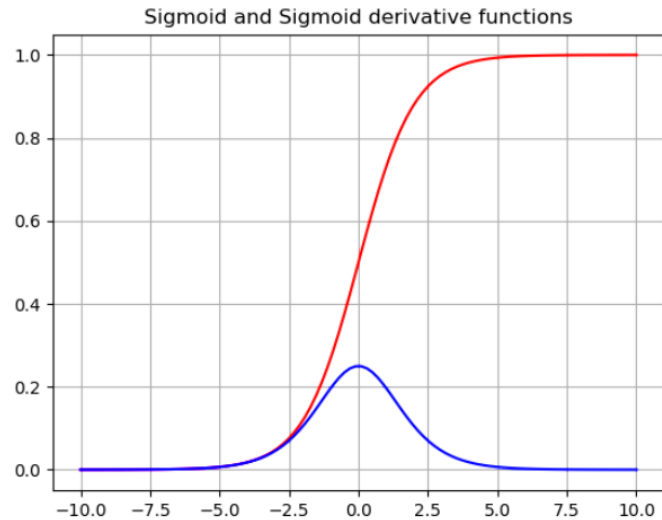
# Optimization Challenges

- › Exploding Gradient in Cliff Structure



# Optimization Challenges

- › In Deep Neural Network:



- › Chaotic Behavior in Nonlinear Function!
- › ReLU is a good choice!
- › Gradient Clipping ( $\alpha \frac{g}{\|g\|}$ )



# Stochastic Gradient Descent (SGD)

---

Algorithm 8.1 Stochastic gradient descent (SGD) update at training iteration  $k$

---

Require: Learning rate  $\epsilon_k$ .

Require: Initial parameter  $\theta$

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

    Compute gradient estimate:  $\hat{g} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

    Apply update:  $\theta \leftarrow \theta - \epsilon \hat{g}$

end while

---

$$\epsilon_k = (1 - \alpha)\epsilon_0 + \alpha\epsilon_{\tau} \quad (8.14)$$

with  $\alpha = \frac{k}{\tau}$ . After iteration  $\tau$ , it is common to leave  $\epsilon$  constant.



# Momentum

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\boldsymbol{\theta}} \left( \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) \right), \quad (8.15)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{v}. \quad (8.16)$$

---

**Algorithm 8.2** Stochastic gradient descent (SGD) with momentum

---

**Require:** Learning rate  $\epsilon$ , momentum parameter  $\alpha$ .

**Require:** Initial parameter  $\boldsymbol{\theta}$ , initial velocity  $\mathbf{v}$ .

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient estimate:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$

    Compute velocity update:  $\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g}$

    Apply update:  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{v}$

end while

---



# Nesterov Momentum

$$v \leftarrow \alpha v - \epsilon \nabla_{\theta} \left[ \frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta + \alpha v), y^{(i)}) \right], \quad (8.21)$$

$$\theta \leftarrow \theta + v, \quad (8.22)$$

---

**Algorithm 8.3** Stochastic gradient descent (SGD) with Nesterov momentum

---

**Require:** Learning rate  $\epsilon$ , momentum parameter  $\alpha$ .

**Require:** Initial parameter  $\theta$ , initial velocity  $v$ .

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding labels  $y^{(i)}$ .

    Apply interim update:  $\tilde{\theta} \leftarrow \theta + \alpha v$

    Compute gradient (at interim point):  $g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$

    Compute velocity update:  $v \leftarrow \alpha v - \epsilon g$

    Apply update:  $\theta \leftarrow \theta + v$

**end while**

---



# Algorithms with Adaptive Learning Rates

› Idea:

$$\theta_{k+1} = \theta_k - \eta H_k^{-1} \nabla J(\theta_k)$$

$\vdots$

$$\theta_{k+1} = \theta_k - \eta B^{-1} \nabla J(\theta_k), \text{ B: Preconditioner}$$

$\vdots$

$$B_k = \text{diag} \left( \sum_{t=1}^k [\nabla J(\theta_t)] [\nabla J(\theta_t)]^T \right)^{0.5}$$

$$\theta_{k+1} = \theta_k - \eta B_k^{-1} \nabla J(\theta_k) \Rightarrow \theta_{k+1} = \theta_k - \frac{\eta}{\sqrt{B_k + \varepsilon}} \odot \nabla J(\theta_k)$$





# Algorithms with Adaptive Learning Rates

## › AdaGrad

---

### Algorithm 8.4 The AdaGrad algorithm

---

Require: Global learning rate  $\epsilon$

Require: Initial parameter  $\theta$

Require: Small constant  $\delta$ , perhaps  $10^{-7}$ , for numerical stability

Initialize gradient accumulation variable  $r = 0$

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

    Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

    Accumulate squared gradient:  $r \leftarrow r + g \odot g$

    Compute update:  $\Delta\theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{r}} \odot g$ . (Division and square root applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta\theta$

end while

---



# Algorithms with Adaptive Learning Rates

## › RMSProp

---

### Algorithm 8.5 The RMSProp algorithm

---

Require: Global learning rate  $\epsilon$ , decay rate  $\rho$ .

Require: Initial parameter  $\theta$

Require: Small constant  $\delta$ , usually  $10^{-6}$ , used to stabilize division by small numbers.

Initialize accumulation variables  $r = 0$

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

    Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

    Accumulate squared gradient:  $r \leftarrow \rho r + (1 - \rho) g \odot g$

    Compute parameter update:  $\Delta\theta = -\frac{\epsilon}{\sqrt{\delta+r}} \odot g$ . ( $\frac{1}{\sqrt{\delta+r}}$  applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta\theta$

end while

---



# Algorithms with Adaptive Learning Rates

## › RMSProp algorithm with Nesterov momentum

---

### Algorithm 8.6 RMSProp algorithm with Nesterov momentum

---

**Require:** Global learning rate  $\epsilon$ , decay rate  $\rho$ , momentum coefficient  $\alpha$ .

**Require:** Initial parameter  $\theta$ , initial velocity  $v$ .

Initialize accumulation variable  $r = 0$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

    Compute interim update:  $\tilde{\theta} \leftarrow \theta + \alpha v$

    Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$

    Accumulate gradient:  $r \leftarrow \rho r + (1 - \rho) g \odot g$

    Compute velocity update:  $v \leftarrow \alpha v - \frac{\epsilon}{\sqrt{r}} \odot g$ . ( $\frac{1}{\sqrt{r}}$  applied element-wise)

    Apply update:  $\theta \leftarrow \theta + v$

**end while**

---



# ADAM (Adaptive Moment Estimation)

---

Algorithm 8.7 The Adam algorithm

---

**Require:** Step size  $\epsilon$  (Suggested default: 0.001)

**Require:** Exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1)$ .  
(Suggested defaults: 0.9 and 0.999 respectively)

**Require:** Small constant  $\delta$  used for numerical stabilization. (Suggested default:  $10^{-8}$ )

**Require:** Initial parameters  $\theta$

Initialize 1st and 2nd moment variables  $s = \mathbf{0}$ ,  $r = \mathbf{0}$

Initialize time step  $t = 0$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

    Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

$t \leftarrow t + 1$

    Update biased first moment estimate:  $s \leftarrow \rho_1 s + (1 - \rho_1)g$

    Update biased second moment estimate:  $r \leftarrow \rho_2 r + (1 - \rho_2)g \odot g$

    Correct bias in first moment:  $\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}$

    Correct bias in second moment:  $\hat{r} \leftarrow \frac{r}{1 - \rho_2^t}$

    Compute update:  $\Delta\theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$  (operations applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta\theta$

**end while**

---



# ADAM Variation

› ADAMax:

$$r_t = \rho r_{t-1} + (1 - \rho) |g_t|^2$$

$\vdots$

$$r_t = \rho^p r_{t-1} + (1 - \rho^p) |g_t|^p$$

$$p \rightarrow \infty$$

$$r_t = \max \{ \rho r_{t-1}, |g_t| \}$$

› Nesterov-accelerated Adaptive Moment Estimation (NADAM)



## 2nd Order Methods

› Newton's Method Idea:

$$J(\theta) \approx J(\theta_0) + (\theta - \theta_0)^\top \nabla_{\theta} J(\theta_0) + \frac{1}{2}(\theta - \theta_0)^\top \mathbf{H}(\theta - \theta_0), \quad (8.26)$$

$$\theta^* = \theta_0 - \mathbf{H}^{-1} \nabla_{\theta} J(\theta_0) \quad (8.27)$$



# Newton Methods for Convex (pdm H) Loss!

---

**Algorithm 8.8** Newton's method with objective  $J(\theta) = \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$ .

---

Require: Initial parameter  $\theta_0$

Require: Training set of  $m$  examples

while stopping criterion not met do

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$

    Compute Hessian:  $\mathbf{H} \leftarrow \frac{1}{m} \nabla_{\theta}^2 \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$

    Compute Hessian inverse:  $\mathbf{H}^{-1}$

    Compute update:  $\Delta\theta = -\mathbf{H}^{-1} \mathbf{g}$

    Apply update:  $\theta = \theta + \Delta\theta$

end while

---



# Newton Methods for Non pdm H

- › Regularization (Positive  $\alpha$ )

$$\theta^* = \theta_0 - [H(f(\theta_0)) + \alpha I]^{-1} \nabla_{\theta} f(\theta_0). \quad (8.28)$$

- › Conjugate Gradient, BFGS,....





## Additional Strategies for SGD

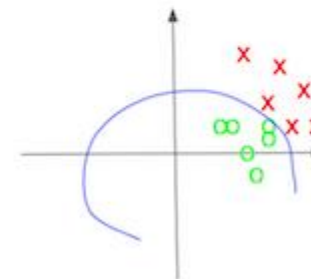
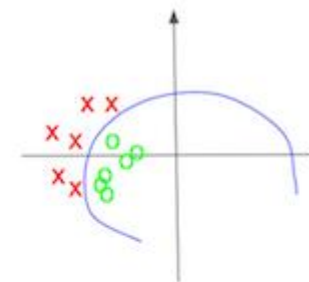
- › Random Shuffle training data after every epoch
- › Curriculum Learning: Sort data (easy to hard)
- › Batch normalization: Re-normalizes every mini-batch to zero mean, unit variance
- › Gradient noise (Decreasing variance)



# Batch Normalization

<https://www.learnopencv.com/batch-normalization-in-deep-networks/>

- › Internal Covariate Shift:
- › Different  $pdf$  in training and test set



# Batch Normalization

<https://www.learnopencv.com/batch-normalization-in-deep-networks/>

- › Need same distribution in each minibatch!
- › How to Solve:
  - **Input layer** normalization (He/Xavier/...) and shuffling!
- › What about hidden layer:
  - Uniform/Normal *dist.* in **input** layer  $\Rightarrow$  Non-uniform/Non-Normal *dist.* In **hidden** layer
  - Hidden layer input *dist.* varied in each batch! (**internal covariate shift**)
- › Original work:
  - *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift* by Sergey Ioffe and Christian Szegedy, 2015 (Google Team)



# Batch Normalization

- › Idea: Insert a (Batch) Normalization Layer (BN-Layer) before each layer
- › Normalization need data → Batch/min or batch training
- › Normalization per cell! (Whitening need too many sample per mini-batch)

Normalize each layer output (input of next layer) in each iteration (min-batch)



# Batch Normalization

- › Training Phase Algorithm:
  - Zero-mean, unit-var is not ideal for nonlinear units
  - Two Learnable extra-weights per cell!
  - Next Layer input is  $y_i$

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots x_m\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$



# Batch Normalization

## › EBP for Two Extra Weights

$$\frac{\partial \ell}{\partial \hat{x}_i} = \frac{\partial \ell}{\partial y_i} \cdot \gamma$$

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}.$$

$$\frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^2} = \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot (x_i - \mu_{\mathcal{B}}) \cdot \frac{-1}{2} (\sigma_{\mathcal{B}}^2 + \epsilon)^{-3/2}$$

$$\frac{\partial \ell}{\partial \mu_{\mathcal{B}}} = \left( \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \right) + \frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^2} \cdot \frac{\sum_{i=1}^m -2(x_i - \mu_{\mathcal{B}})}{m}$$

$$\frac{\partial \ell}{\partial x_i} = \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^2} \cdot \frac{2(x_i - \mu_{\mathcal{B}})}{m} + \frac{\partial \ell}{\partial \mu_{\mathcal{B}}} \cdot \frac{1}{m}$$

$$\frac{\partial \ell}{\partial \gamma} = \sum_{i=1}^m \frac{\partial \ell}{\partial y_i} \cdot \hat{x}_i$$

$$\frac{\partial \ell}{\partial \beta} = \sum_{i=1}^m \frac{\partial \ell}{\partial y_i}$$



# Batch Normalization

## › Test/Inference Phase:

- There is **ONE** sample in input **NOT** mini-batch
- But we need  $E\{x\}$  and  $\text{var}\{x\}$  to normalize

$$\hat{x} = \frac{x - E[x]}{\sqrt{\text{Var}[x] + \epsilon}}$$

1. Use expectation (unbiased estimator) of  $\mu_B$  and  $\sigma_B^2$  instead, as following:

$$\begin{aligned} \mu_B &\leftarrow \frac{1}{m} \sum_{i=1}^m x_i \\ \sigma_B^2 &\leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \end{aligned} \quad \longrightarrow \quad \begin{aligned} E[x] &\leftarrow E_B[\mu_B] \\ \text{Var}[x] &\leftarrow \boxed{\frac{m}{m-1}} E_B[\sigma_B^2] \end{aligned}$$



# Batch Normalization

## 2. Running mean/variance during training phase

$$\hat{\mu}_B(t) = \lambda \hat{\mu}_B(t-1) + (1-\lambda) \mu_B(t-1), \lambda \sim 0.95$$

$$\widehat{\sigma}_B^2(t) = \lambda \widehat{\sigma}_B^2(t-1) + (1-\lambda) \sigma_B^2(t-1)$$

$$y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}^{inf}(x^{(k)}) = \frac{\gamma}{\sqrt{Var[x^{(k)}] + \epsilon}} x^{(k)} + \left( \beta - \frac{\gamma E[x^{(k)}]}{\sqrt{Var[x^{(k)}] + \epsilon}} \right)$$



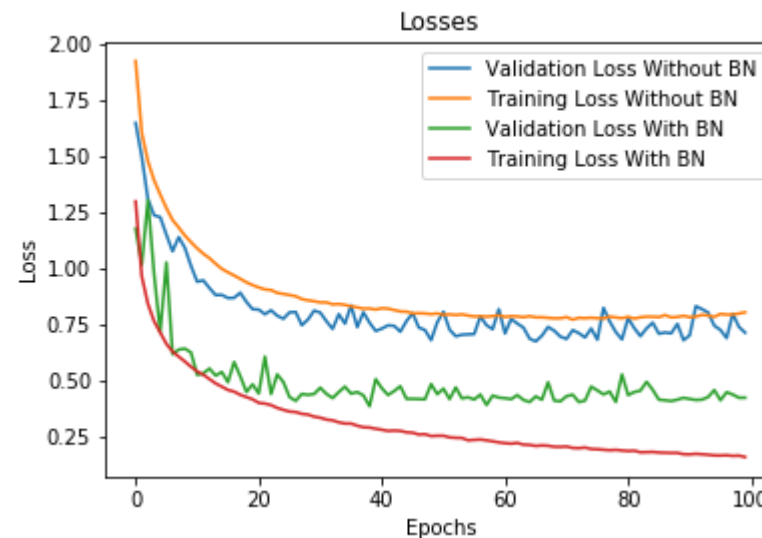


# Batch Normalization

<https://www.learnopencv.com/batch-normalization-in-deep-networks/>

## › Other benefits:

- Higher learning rate (avoiding vanishing/exploding gradient)
- Regularization (Batch information for normalization)
- May use saturating active function



# Regularization

› Any Question?

