

ECE/SIOC 228 Machine Learning for Physical Applications

Lecture 4: Optimization and Regularization

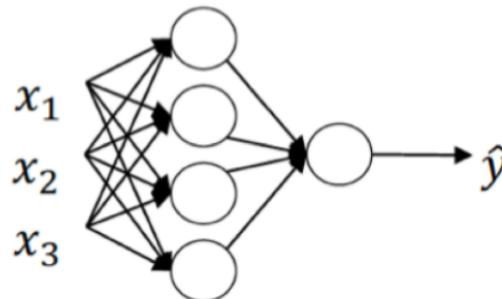
Yuanyuan Shi

Assistant Professor, ECE

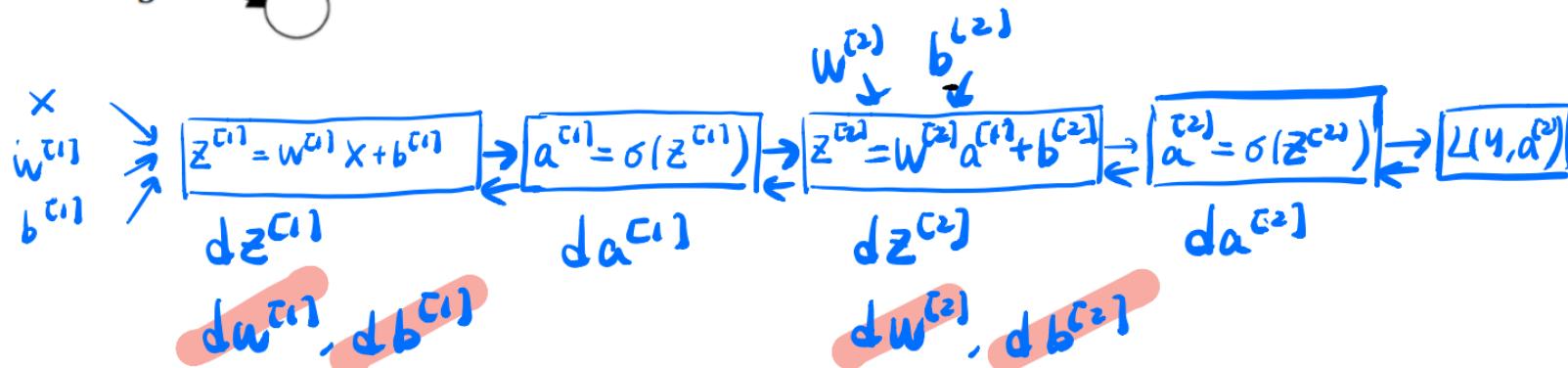
University of California, San Diego

Review: Forward and Backward Computation of Neural Network

UC San Diego



1 sample $(x, y), x \in \mathbb{R}^3, y \in \{0, 1\}$

$$L(\hat{y}, y) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$$


Weight update: $w^{[l]} \leftarrow w^{[l]} - \eta dw^{[l]}, b^{[l]} \leftarrow b^{[l]} - \eta db^{[l]}, l=1, 2$

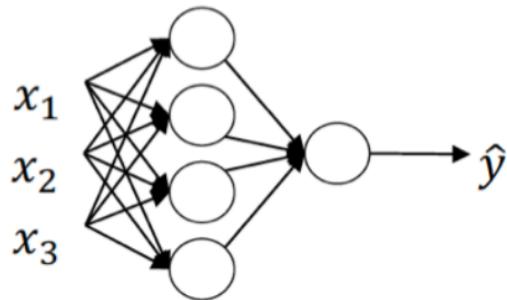
Today's Lecture

UC San Diego

- Forward and Backward Computation across Multiple Examples
- Optimization Algorithms
- Regularization

Vectorizing across Multiple Examples

UC San Diego



for one sample (x, y) , we have:

$$z^{(1)} = W^{(1)}x + b^{(1)}, \quad a^{(1)} = \sigma(z^{(1)})$$

$$z^{(2)} = W^{(2)}a^{(1)} + b^{(2)}, \quad a^{(2)} = \sigma(z^{(2)})$$

What if we have m samples?

$$\begin{aligned} x^{(1)} &\rightarrow \hat{y}^{(1)} \\ x^{(2)} &\rightarrow \hat{y}^{(2)} \\ x^{(m)} &\rightarrow \hat{y}^{(m)} \end{aligned}$$

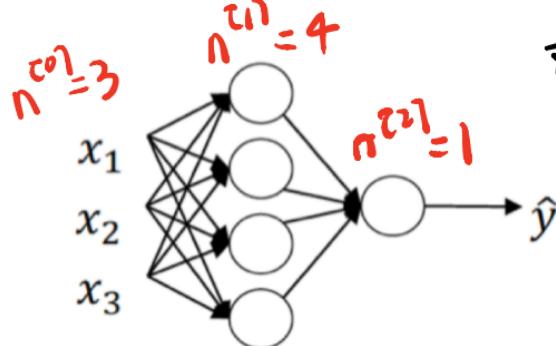
for $i=1$ to m

\downarrow layer 1 \leftarrow sample i

$$z^{(1)(i)} = W^{(1)}x^{(i)} + b^{(1)}, \quad a^{(1)(i)} = \sigma(z^{(1)(i)})$$
$$z^{(2)(i)} = W^{(2)}a^{(1)(i)} + b^{(2)}, \quad a^{(2)(i)} = \sigma(z^{(2)(i)})$$
$$\hat{y}^{(i)} = a^{(2)(i)}$$

Vectorizing across Multiple Examples

UC San Diego



* Speedup computation: Avoid using explicit for-loop (GPU/CPU have parallelizations)

* $\mathbf{z}^{(1)(i)} = \mathbf{w}^{(1)} \mathbf{x}^{(i)} + \mathbf{b}^{(1)}$. Suppose $\mathbf{b}^{(1)}$ is zero vector

$$\mathbf{W}^{(1)} \begin{bmatrix} \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \mathbf{x}^{(3)} & \dots \end{bmatrix} = \begin{bmatrix} \mathbf{w}^{(1)} \mathbf{x}^{(1)} \\ \vdots \\ \mathbf{w}^{(1)} \mathbf{x}^{(2)} \\ \vdots \\ \mathbf{w}^{(1)} \mathbf{x}^{(3)} \\ \vdots \end{bmatrix} = \begin{bmatrix} \mathbf{z}^{(1)(1)} & \mathbf{z}^{(1)(2)} & \mathbf{z}^{(1)(3)} & \dots \end{bmatrix}$$

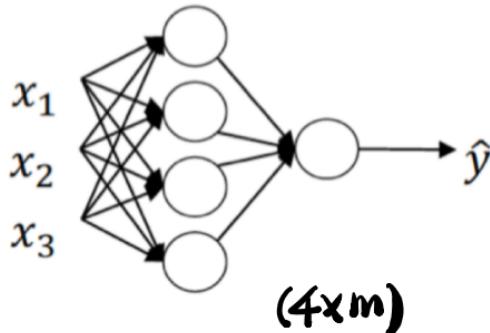
input matrix

$\mathbf{w}^{(1)} \mathbf{x}^{(1)}$

$\mathbf{z}^{(1)}$

Vectorizing across Multiple Examples

UC San Diego



$$\underbrace{W^{(1)}}_{(4 \times 3)} \begin{bmatrix} x^{(1)} \\ | \\ x^{(2)} \\ | \\ x^{(3)} \\ | \\ \dots \end{bmatrix} + \underbrace{\begin{bmatrix} b^{(1)} \\ | \\ b^{(2)} \\ | \\ b^{(3)} \\ | \\ \dots \end{bmatrix}}_{(3 \times m)} \quad \text{(not valid sum!)}$$

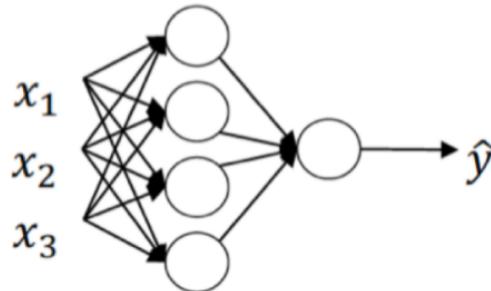
(4×1) ??

$$\Rightarrow \begin{bmatrix} \cdot & \cdot & \cdot \\ \vdots & \vdots & \vdots \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} + \begin{bmatrix} b^{(1)} \\ | \\ b^{(2)} \\ | \\ b^{(3)} \\ | \\ \dots \end{bmatrix} = \begin{bmatrix} z^{(1)(1)} & z^{(1)(2)} & z^{(1)(3)} \\ | & | & | \\ \vdots & \vdots & \vdots \\ \cdot & \cdot & \cdot \end{bmatrix} \quad (4 \times m)$$

automatic
broadcasting in python

Vectorizing across Multiple Examples

UC San Diego



$$X = \begin{bmatrix} X^{(1)} & X^{(2)} & \dots & X^{(m)} \end{bmatrix}$$

$$Z^{(1)} = W^{(1)} X + b^{(1)},$$

$$A^{(1)} = \sigma(Z^{(1)})$$

$$Z^{(2)} = W^{(2)} A^{(1)} + b^{(2)}$$

$$A^{(2)} = \sigma(Z^{(2)})$$

Make sure dimension is right!

$$X : (3 \times m) \quad W^{(1)} : (4 \times 3)$$

$$b^{(1)} : (4 \times 1)$$

$$Z^{(1)} : (4 \times m), A^{(1)} : (4 \times m)$$

$$W^{(2)} : (1 \times 4), b^{(2)} : (1 \times 1)$$

$$Z^{(2)} : (1 \times m), A^{(2)} : (1 \times m)$$

Vectorizing across Multiple Examples

UC San Diego

$$L(\gamma, \hat{\gamma}) = \frac{1}{m} \sum_{i=1}^m L(\gamma^{(i)}, \hat{\gamma}^{(i)}) = \frac{1}{m} \sum_{i=1}^m -[\gamma^{(i)} \log \hat{\gamma}^{(i)} + (1-\gamma^{(i)}) \log (1-\hat{\gamma}^{(i)})]$$

for $i = 1$ to m

$$z^{(1)(i)} = w^{(1)} x^{(i)} + b^{(1)}, a^{(1)(i)} = \sigma(z^{(1)})$$

$$z^{(2)(i)} = w^{(2)} a^{(1)(i)} + b^{(2)}, a^{(2)(i)} = \sigma(z^{(2)})$$

$$\hat{y}^{(i)} = a^{(2)(i)}$$

$$L_t = -[\gamma^{(i)} \log a^{(2)(i)} + (1-\gamma^{(i)}) \log (1-a^{(2)(i)})]$$

$$dz^{(2)(i)} = a^{(2)(i)} - y^{(i)}$$

$$dw^{(2)}_+ = dz^{(2)(i)} a^{(1)(i)T}$$

$$db^{(2)}_+ = dz^{(2)(i)}$$

$$dw^{(1)}_+ = dw^{(2)(i)}, db^{(1)}_+ = db^{(2)(i)}$$

loop over all samples

$$L \leftarrow \frac{1}{m} L$$

$$dw^{(1)} \leftarrow \frac{1}{m} dw^{(1)}$$

$$db^{(1)} \leftarrow \frac{1}{m} db^{(1)}$$

$$dw^{(2)} \leftarrow \frac{1}{m} dw^{(2)}$$

$$db^{(2)} \leftarrow \frac{1}{m} db^{(2)}$$

Vectorizing across Multiple Examples

UC San Diego

$$L(\gamma, \hat{\gamma}) = \frac{1}{m} \sum_{i=1}^m L(\gamma^{(i)}, \hat{\gamma}^{(i)}) = \frac{1}{m} \sum_{i=1}^m -[\gamma^{(i)} \log \hat{\gamma}^{(i)} + (1-\gamma^{(i)}) \log (1-\hat{\gamma}^{(i)})]$$

for $i = 1$ to m

$$z^{(c_1)(i)} = w^{(1)} x^{(i)} + b^{(1)}, \quad a^{(c_1)(i)} = \sigma(z^{(c_1)})$$

$$z^{(c_2)(i)} = w^{(2)} a^{(c_1)(i)} + b^{(2)}, \quad a^{(c_2)(i)} = \sigma(z^{(c_2)})$$

$$\hat{y}^{(i)} = a^{(c_2)(i)}$$

$$L_i = -[\gamma^{(i)} \log a^{(c_2)(i)} + (1-\gamma^{(i)}) \log (1-a^{(c_2)(i)})]$$

- $d z^{(c_2)(i)} = a^{(c_2)(i)} - y^{(i)}$

- $d w^{(2)} + = d z^{(c_2)(i)} a^{(c_1)(i)} T$

- $d b^{(c_2)} + = d z^{(c_2)(i)}$

- $d w^{(c_1)} + = d w^{(c_1)(i)}, \quad d b^{(c_1)} + = d b^{(c_1)(i)}$

$$d z^{(c_2)} = [a^{(c_2)(1)} - y^{(1)}, a^{(c_2)(2)} - y^{(2)}, \dots, a^{(c_2)(m)} - y^{(m)}]$$

$$= A^{(2)} - Y$$

$$d b^{(c_2)} = \frac{1}{m} \sum_{i=1}^m d z^{(c_2)(i)}$$

$$d w^{(c_2)} = \frac{1}{m} \sum_{i=1}^m d z^{(c_2)(i)} a^{(c_1)(i)} T$$

$$(1 \times 4) = \frac{1}{m} [d z^{(c_2)(1)} \cdots d z^{(c_2)(m)}] \begin{bmatrix} -a^{(c_1)(1)} \\ \vdots \\ -a^{(c_1)(m)} \end{bmatrix}$$

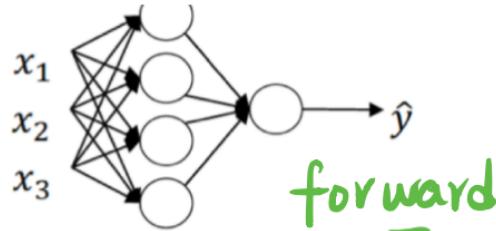
$$= \frac{1}{m} d z^{(c_2)} (A^{(c_1)})^T \begin{bmatrix} -a^{(c_1)(1)} \\ \vdots \\ -a^{(c_1)(m)} \end{bmatrix}$$

$$M \times 4$$

Vectorizing across Multiple Examples

UC San Diego

$$L(\gamma, \hat{\gamma}) = -\frac{1}{m} \sum_{i=1}^m [\gamma^{(i)} \log \hat{\gamma}^{(i)} + (1-\gamma^{(i)}) \log (1-\hat{\gamma}^{(i)})]$$



forward

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \end{bmatrix}$$

$$\mathbf{z}^{(1)} = \mathbf{w}^{(1)} \mathbf{x} + \mathbf{b}^{(1)},$$

$$\mathbf{a}^{(1)} = \sigma(\mathbf{z}^{(1)})$$

$$\mathbf{z}^{(2)} = \mathbf{w}^{(2)} \mathbf{a}^{(1)} + \mathbf{b}^{(2)}$$

$$\mathbf{a}^{(2)} = \sigma(\mathbf{z}^{(2)})$$

Backpropagation

$$d\mathbf{z}^{(2)} = \mathbf{a}^{(2)} - \gamma$$

$$d\mathbf{w}^{(2)} = \frac{1}{m} d\mathbf{z}^{(2)} (\mathbf{a}^{(1)})^T$$

$$db^{(2)} = \frac{1}{m} \sum_{i=1}^m d\mathbf{z}^{(2)(i)}$$

$$d\mathbf{z}^{(1)} = \mathbf{w}^{(2)T} d\mathbf{z}^{(2)} * \sigma'(\mathbf{z}^{(1)})$$

$$d\mathbf{w}^{(1)} = \frac{1}{m} d\mathbf{z}^{(1)} \mathbf{x}^T$$

$$db^{(1)} = \frac{1}{m} \sum_{i=1}^m d\mathbf{z}^{(1)(i)}$$

Mini-batch and Stochastic Gradient Descent

Batch vs. Mini-batch Gradient Descent

UC San Diego

- Vectorization allows you to efficiently compute on m examples

$$X = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & x_1^{(3)} & \dots & x_1^{(m)} \end{bmatrix}$$

What if m
is very large?

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(m)} \end{bmatrix}$$

e.g. m = 1,000,000

Mini-batch

$$X = \underbrace{\begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(1000)} \end{bmatrix}}_{\dots} \begin{bmatrix} x^{(1001)} \\ \vdots \\ x^{(2000)} \end{bmatrix} \dots \begin{bmatrix} \dots \\ x^{(m)} \end{bmatrix}$$
$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(1000)}], [y^{(1001)}, \dots, y^{(2000)}], \dots, [y^{(m)}]$$

Mini-batch Gradient Descent

UC San Diego

repeat } ← epoch

for $t = 1, \dots, 1000$ }

Forward prop on $X^{(t)}$

$$Z^{(1)} = W^{(1)} X^{(t)} + b^{(1)}$$

$$A^{(1)} = g^{(1)}(Z^{(1)})$$

:

$$A^{(L)} = g^{(L)}(Z^{(L)})$$

$$\text{Compute cost } J^{(t)} = \frac{1}{1000} \sum_{i=1}^l L(\hat{y}^{(i)}, y^{(i)})$$

iterate over all mini-batches

backprop to compute gradient using

$$(X^{(t)}, y^{(t)})$$

$$W^{(l)} = W^{(l)} - \eta \nabla W^{(l)}$$

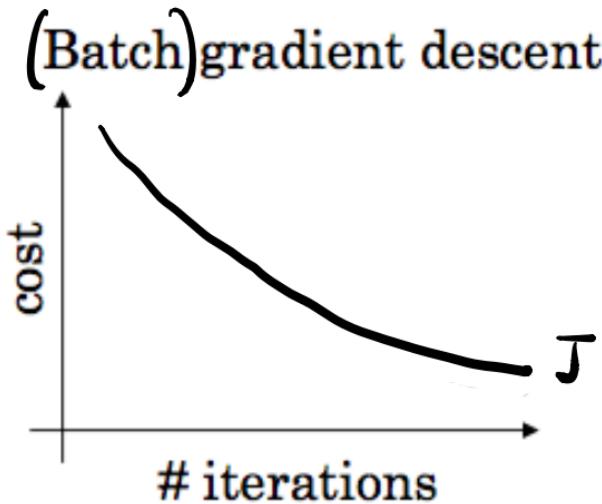
η = learning rate

$$b^{(l)} = b^{(l)} - \eta \nabla b^{(l)}$$

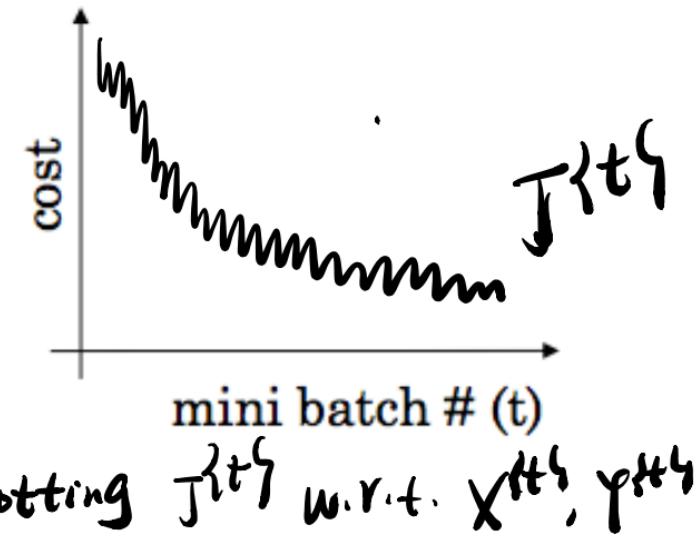
} ← epoch

Understanding Mini-batch Gradient Descent

UC San Diego



Mini-batch gradient descent



Choosing Mini-batch Size

UC San Diego

- If mini-batch size = m : (Batch) gradient descent
- If mini-batch size = 1: stochastic gradient descent
- If mini-batch size $1 < b < m$: mini-batch gradient descent

SGD	Mini-batch gradient descent	Batch gradient descent
	typical size: 64, 128, 256 ... (power of 2)	

$$\min_w L(w) = \frac{1}{m} \sum_{i=1}^m L_i(w)$$

$L_i(w)$ is the loss on sample $(x^{(i)}, y^{(i)})$

Gradient descent will do:

$$w \leftarrow w - \eta \cdot \frac{1}{m} \sum_{i=1}^m \nabla L_i(w)$$

SGD : $w \leftarrow w - \eta_k \cdot \nabla L_{i_k}(w)$
 $i_k \in \{1, 2, \dots, m\}$

* Randomized rule is common in practice;
choose $i \in \{1, 2, \dots, m\}$ uniformly at random

$$E[\nabla L_{ik}(w)] = \nabla L(w)$$

- SGD is using an unbiased estimate of the gradient at each step (Gradient descent + ^{zero-mean} noise)
- Big saving in memory & computation

Mini-batch gradient descent :

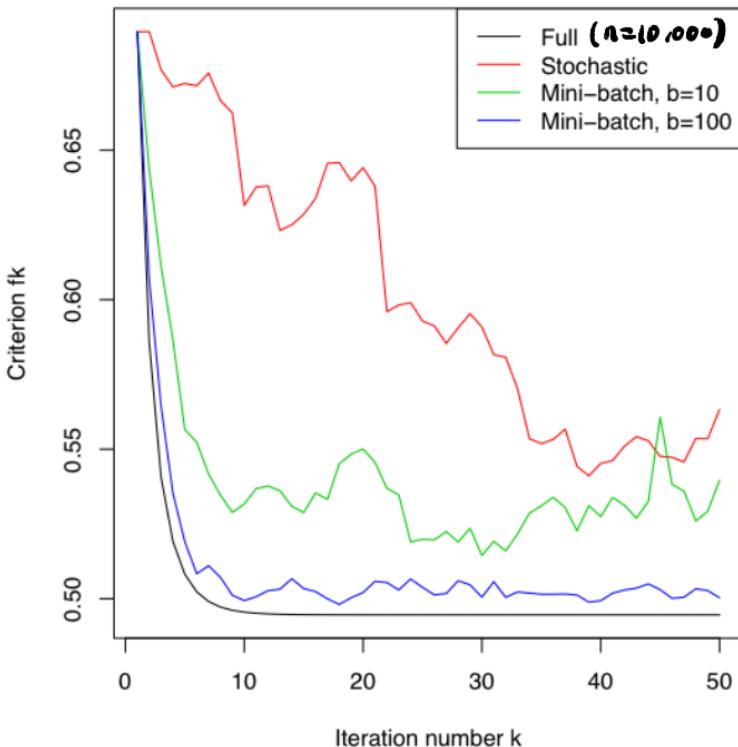
$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{1}{b} \sum_{i \in I_k} \nabla L_i(\mathbf{w})$$

- $I_k \subseteq \{1, \dots, m\}$ is a random subset from the whole dataset. $|I_k| = b \ll m$

- $E\left[\frac{1}{b} \sum_{i \in I_k} \nabla L_i(\mathbf{w})\right] = \nabla L(\mathbf{w})$
 - unbiased
 - reduce variance

SGD Convergence

UC San Diego



- Logistic regression
 - $\hat{y} = \sigma(w^T x + b)$, $x \in \mathbb{R}^{20}$
 - $n = 10,000$ data points
- * SGD/mini-batch gradient descent is the dominant algorithm in training neural networks!

Today's Lecture

UC San Diego

- Forward and Backward Computation across Multiple Examples
- Optimization Algorithms
- Regularization

Optimization Algorithms

Gradient Descent with Momentum

UC San Diego



momentum (previous gradient information)

On iteration t :

- Compute dw, db on the current mini-batch

- $V_{dw} = \beta V_{dw} + (1-\beta)dw$

$$V_{db} = \beta V_{db} + (1-\beta)db$$

momentum

η : learning rate

$$w := w - \eta V_{dw}$$

$$b := b - \eta V_{db}$$

* Usually trains faster

than standard gradient descent

Gradient Descent with Momentum

UC San Diego

Exponential averaging:

$$V_t = \beta V_{t-1} + (1-\beta) g_t$$

$\beta = 0.9$ ← hyperparameter

$$V_{10} = 0.9 V_9 + 0.1 g_{10}$$

$$V_9 = 0.9 V_8 + 0.1 g_9$$

$$V_8 = 0.9 V_7 + 0.1 g_8$$

$$V_{10} = 0.1 g_{10} + 0.9 (0.9 V_8 + 0.1 g_9)$$

$$= 0.1 g_{10} + 0.09 g_9 + 0.81 V_8 (0.1 g_8 + 0.9 V_7)$$

$$= \underbrace{0.1 g_{10}}_{\text{bias term}} + \underbrace{0.1 \times 0.9 g_9}_{\text{dampening term}} + \underbrace{0.1 \times (0.9)^2 g_8}_{\text{dampening term}} + \underbrace{0.1 \times (0.9)^3 g_7}_{\text{dampening term}} + \dots$$

* momentum is an exponential moving average of past gradients

Gradient Descent with Momentum

UC San Diego

Exponential averaging:

$$V_t = \beta V_{t-1} + (1-\beta) g_t$$

← is it the "right" exponential average of past gradient?

$$V_0 = 0$$

$$V_1 = 0.9 V_0 + \underline{0.1 g_1} = 0.1 g_1$$

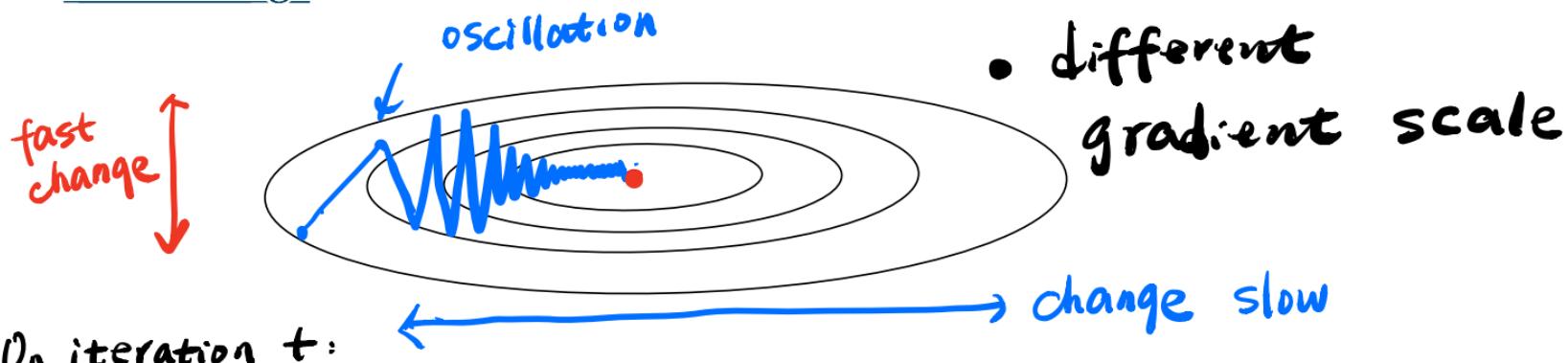
$$V_2 = 0.81 V_0 + \underline{0.1 g_2} + \underline{0.09 g_1} = 0.1 g_2 + 0.09 g_1$$

$$V_t = \beta^t V_0 + \underline{(1-\beta) g_t} + \underline{(1-\beta)\beta g_{t-1}} + \underline{(1-\beta)\beta^2 g_{t-2}} + \dots$$

Sum of Coffs: $1 - \beta^t \Rightarrow V_t^{\text{corrected}} = \frac{V_t}{1 - \beta^t}$

RMSprop (Root Mean Squared propagation)

UC San Diego



On iteration t :

- Compute $d\mathbf{w}$, $d\mathbf{b}$ on the current mini-batch

$$S_{d\mathbf{w}} = \beta_2 S_{d\mathbf{w}} + (1 - \beta_2) d\mathbf{w}^2$$

$$S_{d\mathbf{b}} = \beta_2 S_{d\mathbf{b}} + (1 - \beta_2) d\mathbf{b}^2$$

Adagrad is very similar:

$$S_{d\mathbf{w}} = S_{d\mathbf{w}} + (d\mathbf{w})^2$$

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{d\mathbf{w}}{\sqrt{S_{d\mathbf{w}} + \epsilon}}, \quad \mathbf{b} \leftarrow \mathbf{b} - \eta \frac{d\mathbf{b}}{\sqrt{S_{d\mathbf{b}} + \epsilon}}, \quad \epsilon \text{ is a small constant } (10^{-8})$$

Adam Optimizer

UC San Diego

★ Combines the best properties of "momentum" and "Rmsprop"

V_{dw} = 0, V_{db} = 0, S_{dw} = 0, S_{db} = 0 root mean squared term
momentum term On iteration t:

Compute dw, db using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw \quad \text{"momentum"}$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2 \quad \text{"RmsProp"}$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2$$

Adam Optimizer

UC San Diego

$$w \leftarrow w - \eta \frac{V_{dw}}{\sqrt{S_{dw}} + \epsilon}, \quad b \leftarrow b - \eta \frac{V_{db}}{\sqrt{S_{db}} + \epsilon}$$

+ bias correction

$$(V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1 t))$$

$$(V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1 t))$$

$$(S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2 t))$$

$$(S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2 t))$$

η : needs to be tune

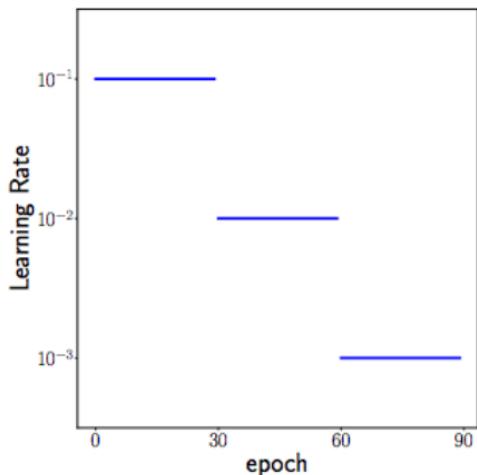
$$\beta_1: 0.9 (dw)$$

$$\beta_2: 0.999 (dw^2)$$

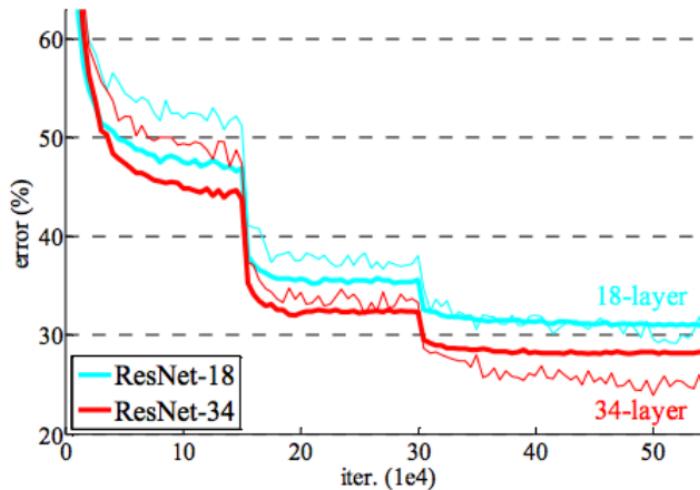
$$\epsilon: 10^{-8}$$

Learning Rate Decay

UC San Diego



(a) Learning rate decay strategy



(b) Figure taken from He et al. (2016)

Figure 1: Training error in (b) is shown by thin curves, while test error in (b) by bold curves.

[1] You, Kaichao, Mingsheng Long, Jianmin Wang, and Michael I. Jordan. "How does learning rate decay help modern neural networks?." *arXiv preprint arXiv:1908.01878* (2019).

Learning Rate Decay

UCSanDiego

- Staircase (reduce by certain factor every few epochs)
- Exponential decaying : $\eta = \eta_0 * \text{decay}^{\text{epoch } \#}$
- Epoch number based:
 - $\eta = \frac{1}{1 + \text{decay} * \text{epoch } \#} * \eta_0$
 - $\eta = \eta_0 \frac{C}{\sqrt{\text{epoch-num}}}$
- Mini-batch number based : $\eta = \eta_0 \frac{C}{\sqrt{t}}$
- Manual decay

Problem of Local Optima

Landscape Conjecture

UC San Diego

- For many years, most practitioners believed that local minima were a common problem in training deep neural networks
- Today, that doesn't seem to be the case – remain an active research area
- Researchers suspect that: for large-size neural networks, **most local minima have a low cost function value and are equivalent**
- Thus, finding the true global minimum on the training set (as opposed to one of the many good local ones) is not very important and may lead to overfitting

Further Reading:

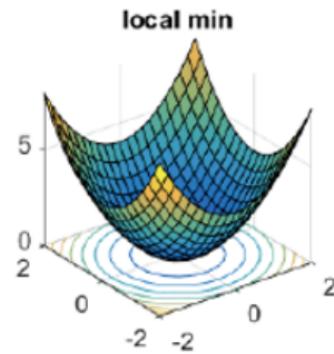
The Loss Surfaces of Multilayer Networks, <https://arxiv.org/pdf/1412.0233.pdf>

<https://www.deeplearningbook.org/contents/optimization.html>

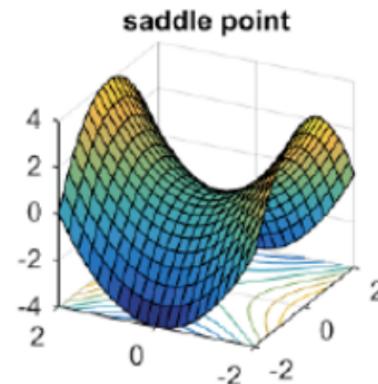
Problem of Saddle Points

UC San Diego

- For high-dimensional nonconvex optimization, local minima are in fact rare compared to another kind of point with zero gradient: saddle points



Hessian: positive definite



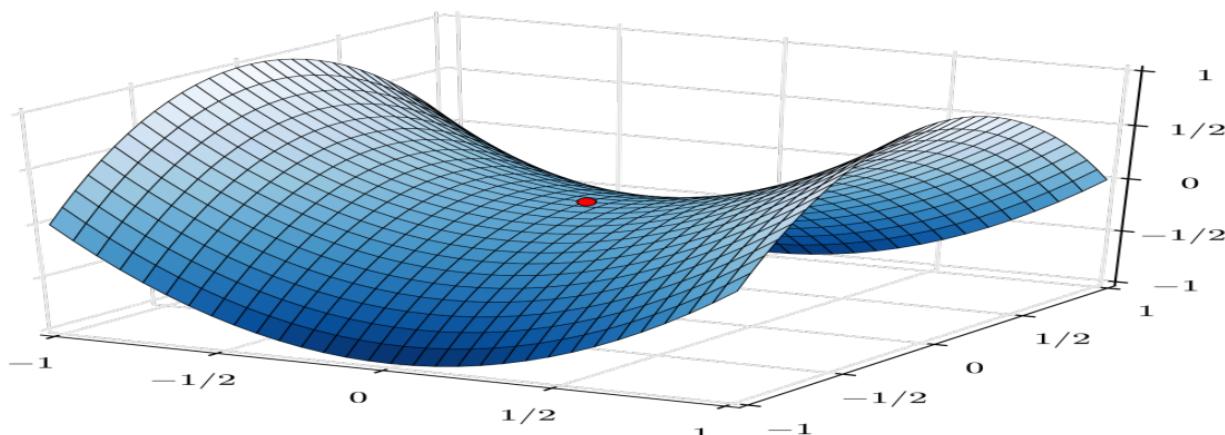
Hessian has both positive/negative eigenvalues

Further Reading:

How to Escape Saddle Points Efficiently <https://arxiv.org/abs/1703.00887>
<https://www.deeplearningbook.org/contents/optimization.html>

Problem of Plateaus

UC San Diego



- Plateaus can slow down learning

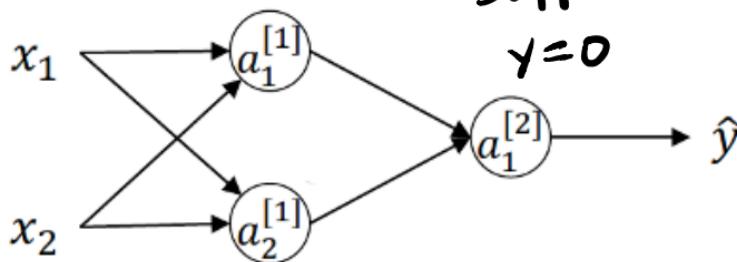
Random Initialization

UC San Diego

$$x \in \mathbb{R}^2 \quad z^{[1]} = W^{[1]} x + b^{[1]}$$

Suppose

$$y=0$$



$$z^{[1]} = W^{[1]} x + b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a^{[1]} = \text{sigmoid}(z^{[1]}) = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} = 0$$

$$a^{[2]} = \text{sigmoid}(z^{[2]}) = \frac{1}{2}$$

What if we initialize all weights to 0?

$$\underline{W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}}, \underline{b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}}$$

$$\underline{W^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix}}, \underline{b^{[2]} = 0}$$

$$\underline{dZ^{[2]} = a^{[2]} - y = \frac{1}{2}}$$

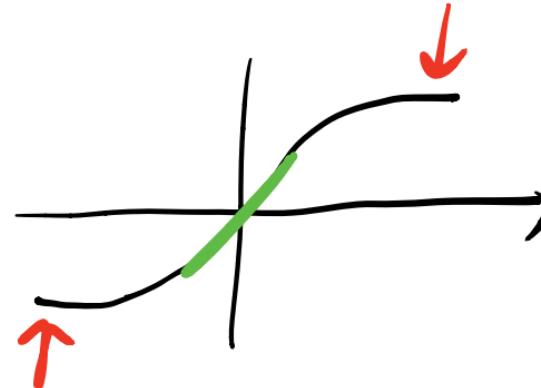
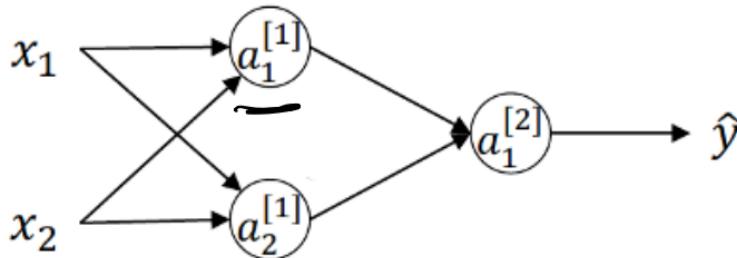
$$\underline{dW^{[2]} = dZ^{[2]} \cdot a^{[1]T}}$$

same for all weights!

$$= \frac{1}{2} \cdot \left[\frac{1}{2} \ \frac{1}{2} \right] = \underline{\underline{\begin{bmatrix} \frac{1}{4} & \frac{1}{4} \end{bmatrix}}}$$

Random Initialization

UC San Diego



- If weights are initialized too large or too small network won't learn well
- Xavier initialization, $\left(\pm \frac{\sqrt{6}}{\sqrt{n_{in} + n_{out}}}\right)$; Kaiming initializer

Today's Lecture

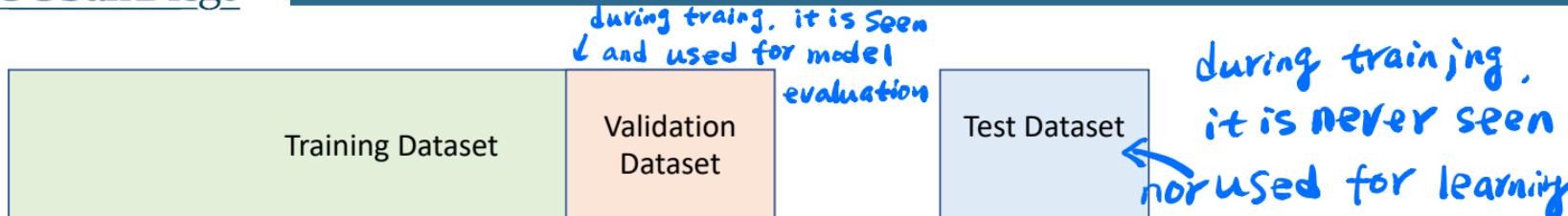
UC San Diego

- Forward and Backward Computation across Multiple Examples
- Optimization Algorithms
- **Regularization**

Regularization

Train, Validation and Test Sets

UC San Diego



[Classic: 70/15/15, 60/20/20] [Big data: 10,000 validators → test] remaining-training

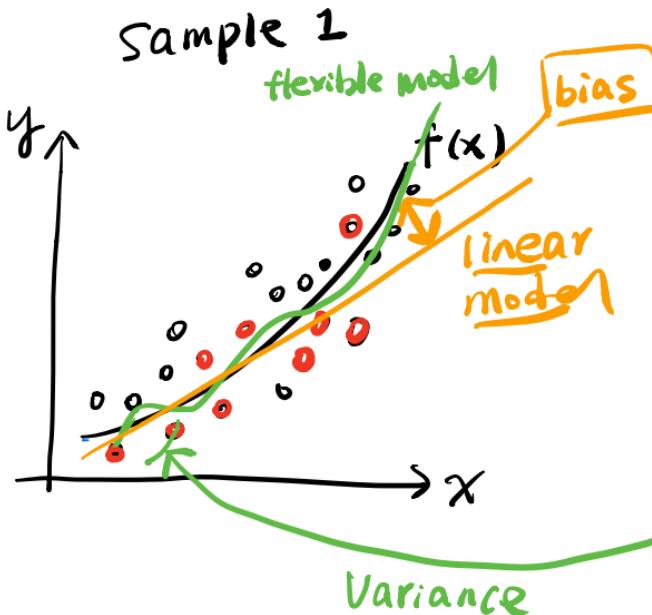
- **Training data**: dataset we use to train the model (e.g. learning weights and biases in a NN model)
- **Validation data (also called development data)**: dataset we use to evaluate model during training & model hyperparameter tuning
- **Test data**: dataset to provide evaluation of the final model fit

Bias and Variance Tradeoff

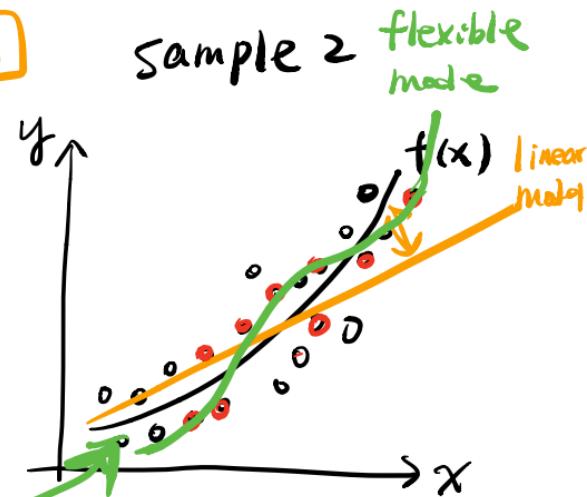
UC San Diego



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



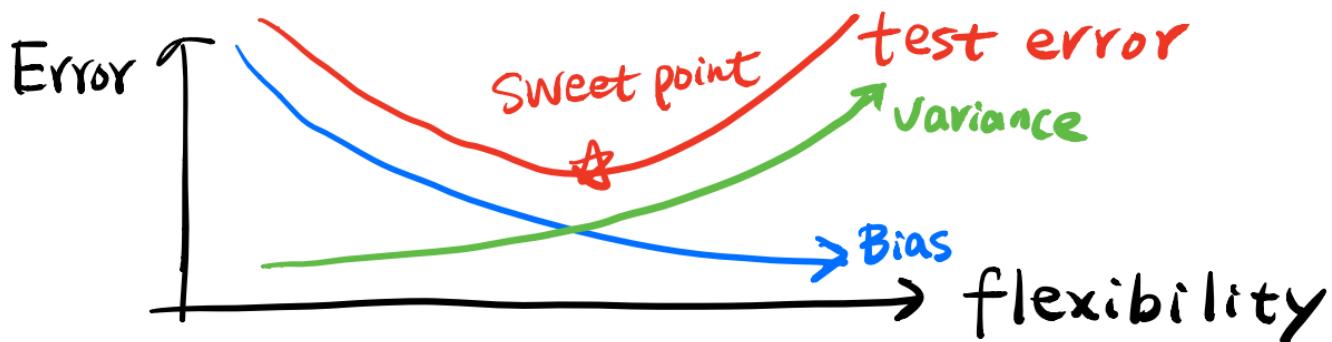
ϵ : zero-mean random noise



Bias and Variance Tradeoff

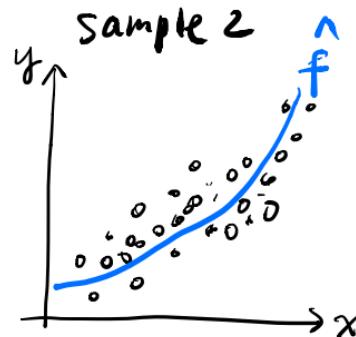
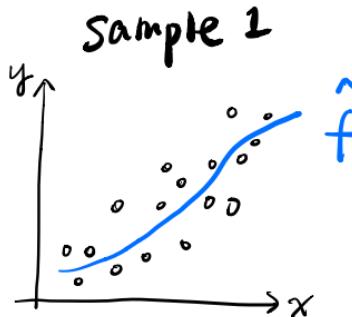
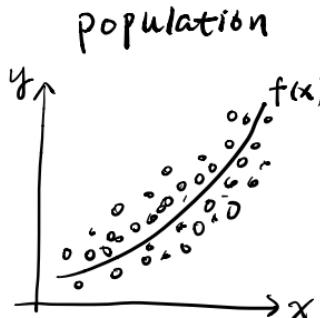
UC San Diego

- **Variance** refers to the amount by which \hat{f} would change if we estimated it using a different training data set. In general, more flexible statistical methods have higher variance.
- **Bias** refers to the error that is introduced by approximating the data (which may arise in a real-life problem with very complicated relationship), by a much simpler model.
- As a general rule, as we use more flexible methods, **the variance will increase and the bias will decrease**.



Bias and Variance Tradeoff

UC San Diego



- For a given datapoint (x, y) underlying relationship as $y = f(x) + \epsilon$, ϵ is a zero mean random noise. The expected test MSE of a learned model \hat{f} can always be decomposed into the sum of three fundamental quantities:

$$E(y - \hat{f}(x))^2 = E(y - E[\hat{f}(x)] + E[\hat{f}(x)] - \hat{f}(x))^2 = E\left(\left(E[\hat{f}(x)] - \hat{f}(x)\right)^2\right) + E\left(\left(f(x) + \epsilon - E[\hat{f}(x)]\right)^2\right)$$

$$= E\left(\left(E[\hat{f}(x)] - \hat{f}(x)\right)^2\right) + E\left(\left(f(x) - E[\hat{f}(x)]\right)^2\right) + Var(\epsilon^2)$$

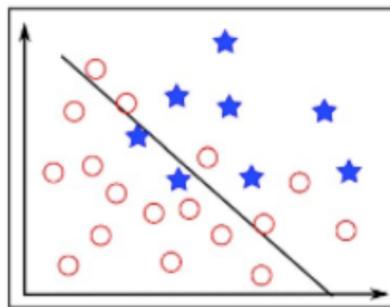
- irreducible

Variance bias

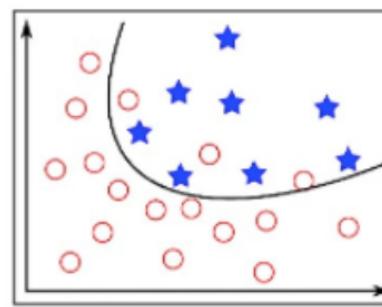
- In order to minimize the expected test error, we need to select a statistical learning method that simultaneously achieves *low variance* and *low bias*.

Bias and Variance

UC San Diego



Underfitting



Overfitting

High bias

High variance

Bias and Variance Tradeoff

UC San Diego

	Error	High bias	Low variance Good model!
Training dataset	1%	14%	1.3%
Validation data	15%	15%	1.3%

High variance
(overfitting)

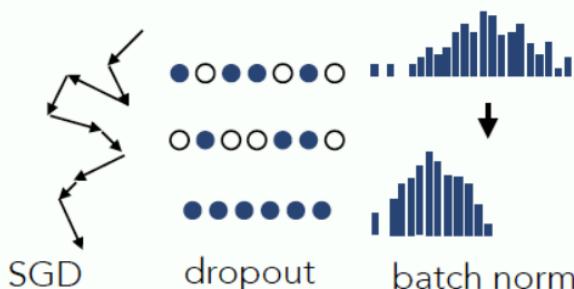
Regularization Methods in Deep Learning

UC San Diego

neural networks are amazingly **flexible**...
given enough parameters, they can perfectly fit random noise

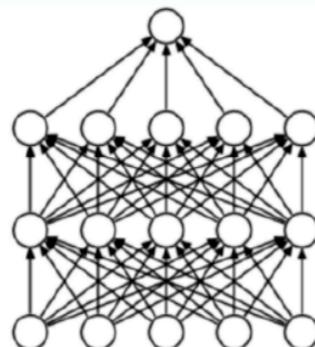
regularization combats **overfitting**

stochasticity (uncertainty)

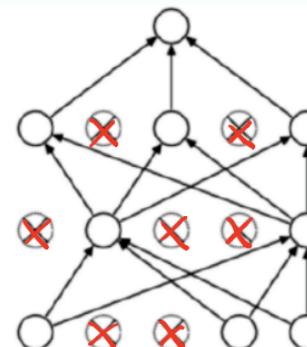


Dropout

UC San Diego



standard neural net



neural net with dropout

0.5

Dropping out units (both hidden and observed) in a neural network

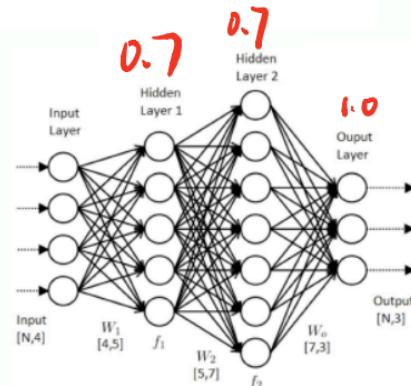
Keep with probability p , drop out with probability $1-p$

Why Does it Work?

UC San Diego

Working with a reduced network at each iteration

Can't rely on one features, spread out the weights



Use varying keep probability for different layers

Wager, Stefan, Sida Wang, and Percy S. Liang. "Dropout training as adaptive regularization." *Advances in neural information processing systems*. 2013.

Batch Normalization

UC San Diego

- Batch norm · normalizes each layer's activations according to statistics of the batch

$$a^{[l-1]} \xrightarrow{w^{[l]}, b^{[l]}} z^{[l]} \longrightarrow \tilde{z}^{[l]} \xrightarrow{\gamma, \beta} a^{[l]}$$

$$\mu_B^{[l]} = \frac{1}{b} \sum_{i \in B} z^{[l](i)}$$

$$\sigma_B^{[l]2} = \frac{1}{b} \sum_{i \in B} (z^{[l](i)} - \mu_B^{[l]})^2$$

(applied to each dim of $z^{[l]}$, independently)

• Normalize :

$$\frac{z^{[l]} - \mu_B^{[l]}}{\sqrt{\sigma_B^{[l]2} + \epsilon}}$$

• Scale &

$$\gamma * \frac{z^{[l]} - \mu_B^{[l]}}{\sqrt{\sigma_B^{[l]2} + \epsilon}} + \beta$$

Batch Normalization

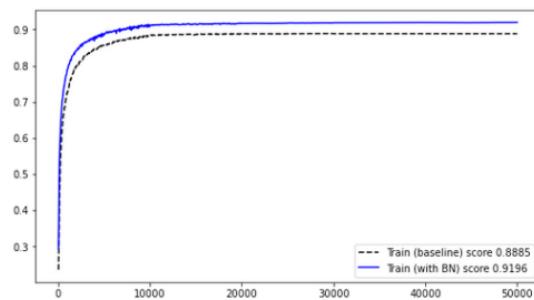
UCSanDiego

- During test: we may not have a full batch to compute the μ_B, σ_B
 - Use $\mu_{\text{pop}}, \sigma_{\text{pop}}$ (estimated μ, σ for the entire training dataset)
- ★ Why batch norm have regularization effect?
use μ, σ from mini-batch (noise)

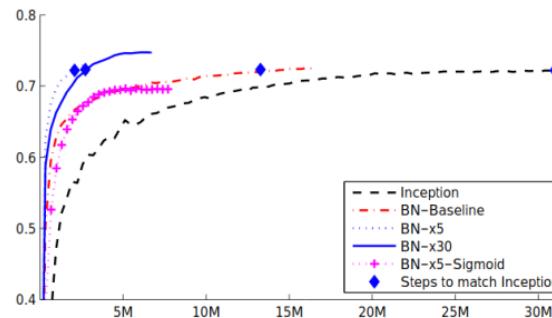
Batch Normalization Helps with Optimization

UC San Diego

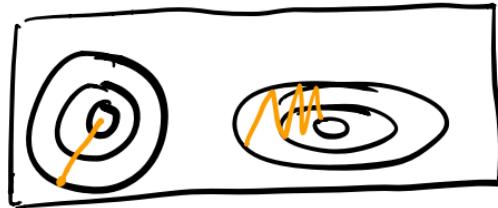
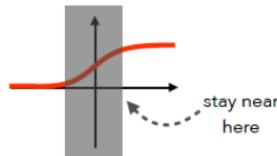
- Batch Normalization is not only a regularization tool but also helps with optimization (**perhaps more important**)



Training 3 fully connected layers of 100 neurons each, all activated by sigmoid function on MNIST



Inception Network on ImageNet



Batch Normalization, Szegedy & Ioffe, 2015

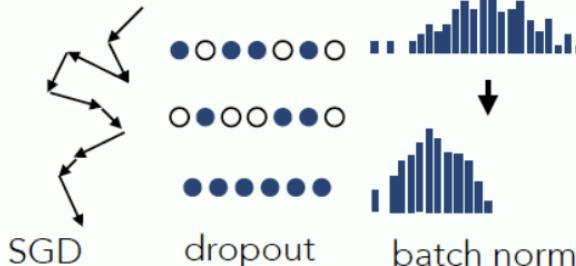
Regularization Methods in Deep Learning

UC San Diego

*neural networks are amazingly flexible...
given enough parameters, they can perfectly fit random noise*

regularization combats **overfitting**

stochasticity (uncertainty)



constraints

