



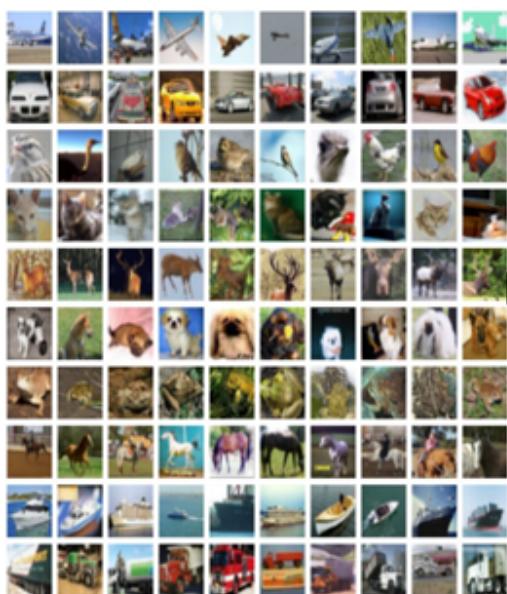
Techniques in Deep Learning

Sourya Dey

Outline

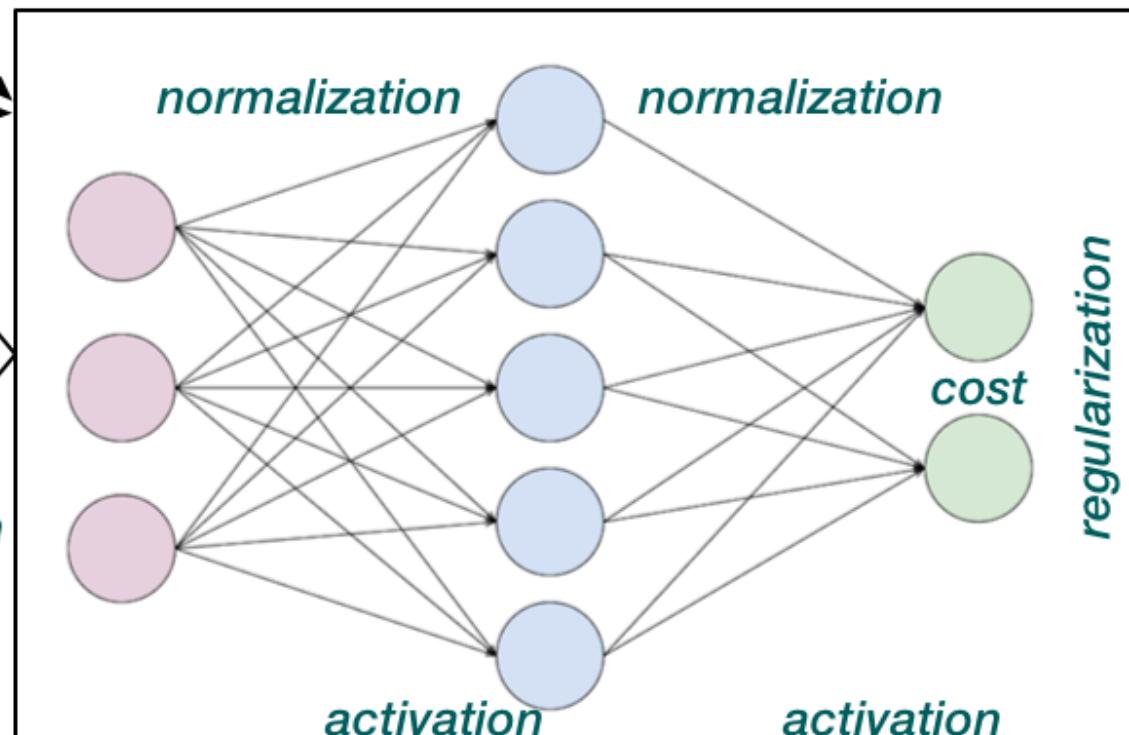
- Activation functions
- Cost function
- Optimizers
- Regularization
- Parameter initialization
- Normalization
- Data handling
- Hyperparameter selection

The Big Picture



*data handling
normalization*

*parameter
initialization*



hyperparameter selection

Outline

- Activation functions
- Cost function
- Optimizers
- Regularization
- Parameter initialization
- Normalization
- Data handling
- Hyperparameter selection

Activation Functions

Recall:

$$s^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}$$

Linear

$$a^{(l)} = h(s^{(l)})$$

Non-Linearity



Non-linearity is required to approximate any arbitrary function

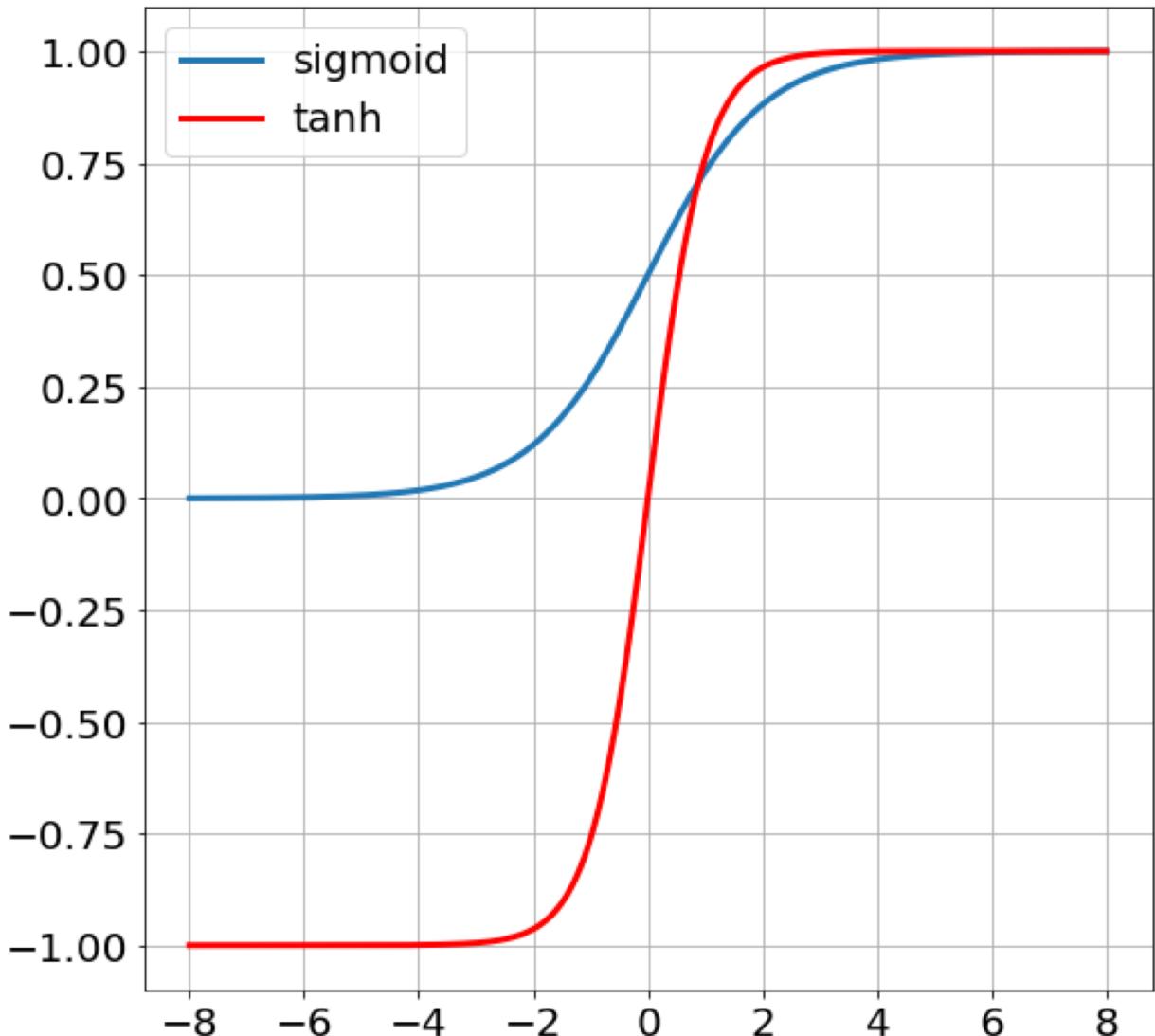
Squashing activations

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

Sigmoid

$$\begin{aligned}\tanh(x) &= \frac{e^x - e^{-x}}{e^x + e^{-x}} \\ &= 2\sigma(2x) - 1\end{aligned}$$

Hyperbolic Tangent
(just rescaled sigmoid)



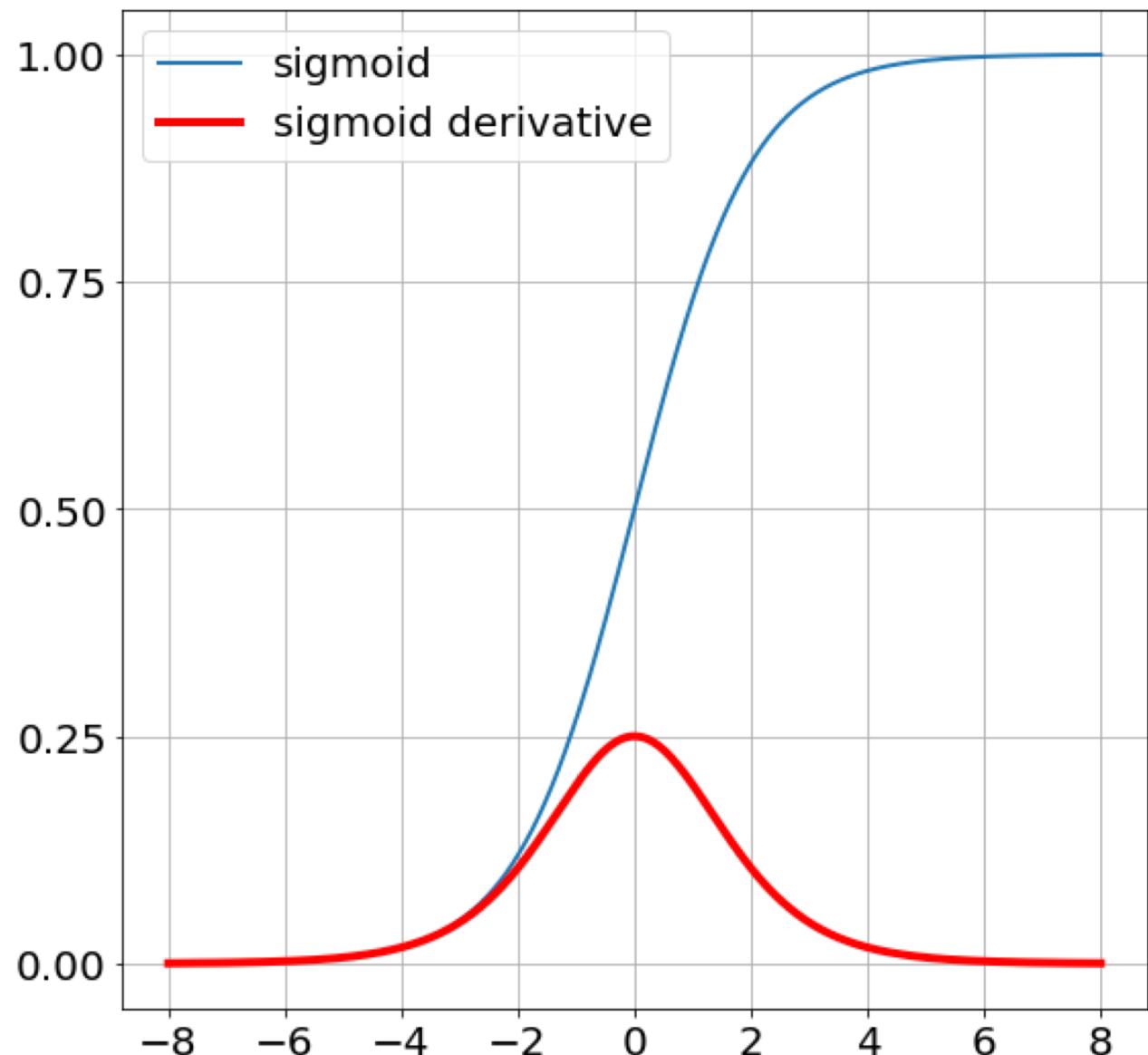
Vanishing gradients

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

Recall from BP (take example L=3)

$$\delta^{(1)} = \mathbf{H}'^{(1)} \mathbf{W}^{(2)T} \mathbf{H}'^{(2)} \mathbf{W}^{(3)T} \mathbf{H}'^{(3)} \nabla_{\mathbf{a}^{(3)}} C$$

All these numbers are <= 0.25, so gradients become very small!!



ReLU family of activations

	$x \geq 0$	$x < 0$
Rectified Linear Unit (ReLU)	x	0
Exponential Linear Unit (ELU)	x	$\alpha(e^x - 1)$
Leaky ReLU	x	αx

Biologically inspired - *neurons firing vs not firing*

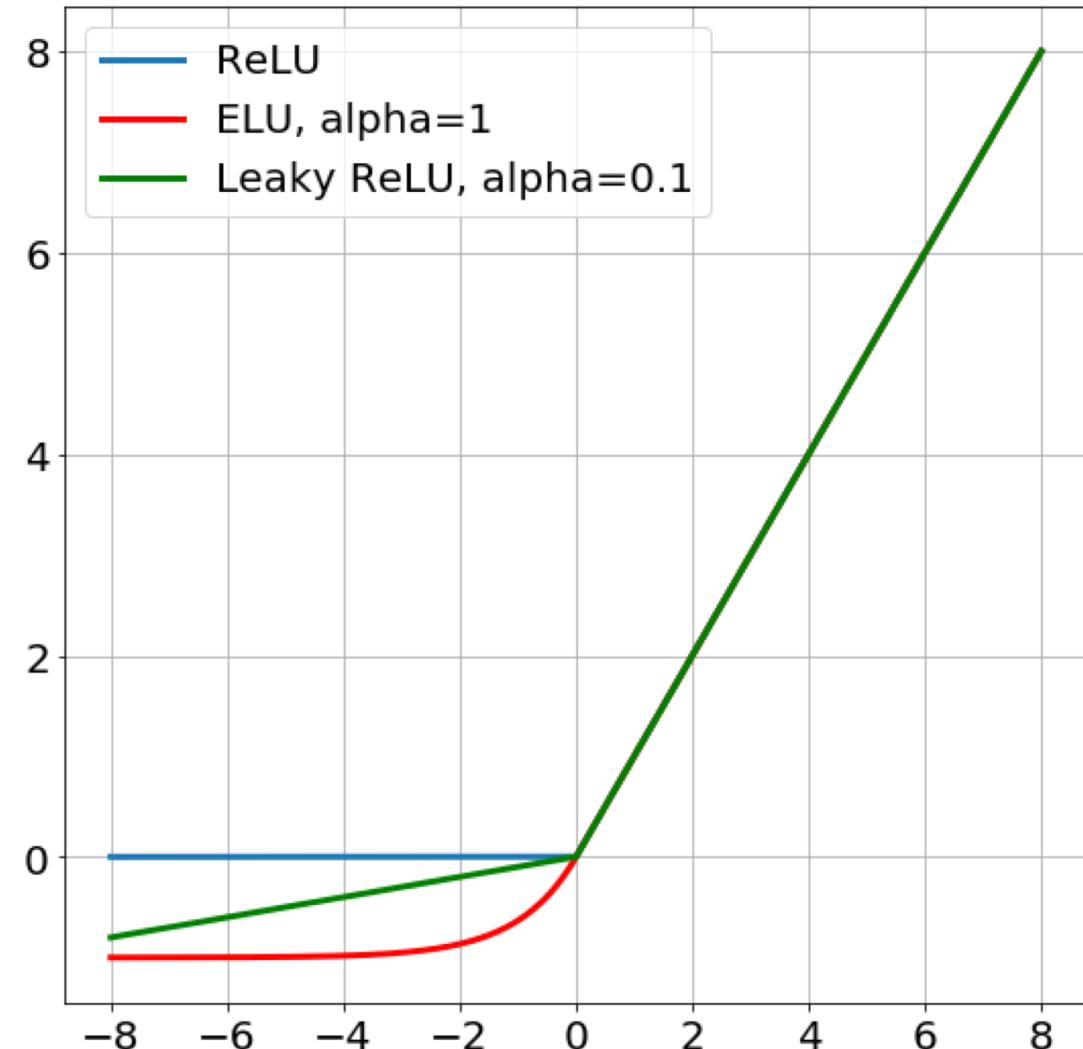
Solves vanishing gradient problem

Non-differentiable at 0, replace with anything in [0,1]

ReLU can die if $x < 0$

Leaky ReLU solves this, but inconsistent results

ELU saturates for $x < 0$, so less resistant to noise



Maxout networks - Generalization of ReLU

Normally:

$$\mathbf{s}^{(l)}_{N^{(l)} \times 1} = \mathbf{W}^{(l)}_{N^{(l)} \times N^{(l-1)}} \mathbf{a}^{(l-1)}_{N^{(l-1)} \times 1} + \mathbf{b}^{(l)}_{N^{(l)} \times 1} \quad \mathbf{a}^{(l)}_{N^{(l)} \times 1} = h\left(\mathbf{s}^{(l)}_{N^{(l)} \times 1}\right)$$

For maxout:

$$\mathbf{s}^{(l)}_{k \times N^{(l)} \times 1} = \mathbf{W}^{(l)}_{k \times N^{(l)} \times N^{(l-1)}} \mathbf{a}^{(l-1)}_{N^{(l-1)} \times 1} + \mathbf{b}^{(l)}_{k \times N^{(l)} \times 1}$$

$$\mathbf{a}^{(l)}_{N^{(l)} \times 1} = \max_{\text{rows}} \mathbf{s}^{(l)}_{k \times N^{(l)} \times 1}$$

Learns the activation function itself

Better approximation power

Takes more computation

Example of maxout

$$k = 2, N^{(l)} = 5$$

$$\begin{matrix} s^{(l)} \\ \scriptstyle 2 \times 5 \end{matrix} = \begin{bmatrix} 2.2 & -1.4 & 9.7 & -1.3 & 0 \\ -4 & -2.1 & 3.8 & -3 & 0.2 \end{bmatrix}$$
$$a^{(l)} = [2.2 \quad -1.4 \quad 9.7 \quad -1.3 \quad 0.2]^T$$

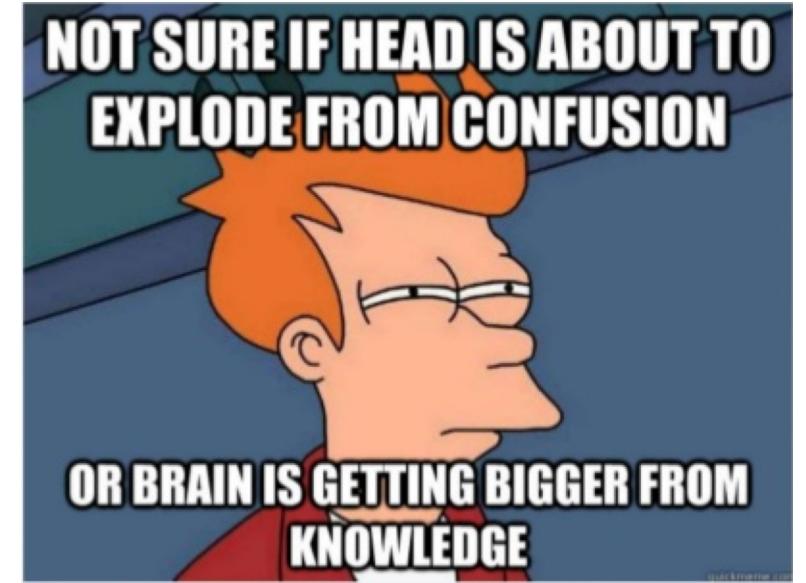
ReLU is a special case of maxout with k=2 and 1 set of (\mathbf{W}, \mathbf{b}) = all 0s

Which activation to use?

Don't use sigmoid

Use ReLU

If too many units are dead, try other activations



And watch out for new activations (or maybe invent a new one) - deep learning moves fast!

Output layer activation - Softmax

$$a = \begin{bmatrix} \frac{e^{s_1}}{\sum_{i=1}^N e^{s_i}} \\ \vdots \\ \frac{e^{s_N}}{\sum_{i=1}^N e^{s_i}} \end{bmatrix}$$

Extending logistic regression to multiple classes

Network output is a probability distribution!

Compare to ideal output probability distribution

$$\mathbf{y}^{(L)} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Outline

- Activation functions
- Cost function
- Optimizers
- Regularization
- Parameter initialization
- Normalization
- Data handling
- Hyperparameter selection

Cross-entropy Cost

$$C = - \sum_{i=1}^{N^{(L)}} y_i^{(L)} \ln a_i^{(L)}$$

Ground truth labels **Network outputs**



Minimizing cross-entropy is the same as minimizing KL-divergence between the probability distributions \mathbf{y} and \mathbf{a}

For binary labels, this reduces to:

$$C = - [y^{(L)} \ln (a^{(L)}) + (1 - y^{(L)}) \ln (1 - a^{(L)})]$$

The one-hot case

$$y^{(L)} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$$

Class 0 incorrect
Class 1 incorrect
Class r correct
Class $N^{(L)}-1$ incorrect

The correct class r is 1, everything else is 0

$$C = -\ln a_r^{(L)}$$

Softmax and cross-entropy with one-hot labels

Recall:

$$\delta^{(L)} = \mathbf{H}'^{(L)T} \nabla_{\mathbf{a}^{(L)}} C$$

$$\delta_r^{(L)} = a_r^{(L)} - 1$$

$$\delta_{\neq r}^{(L)} = a_{\neq r}^{(L)}$$

Combining:

$$\boxed{\delta^{(L)} = \mathbf{a}^{(L)} - \mathbf{y}^{(L)}}$$

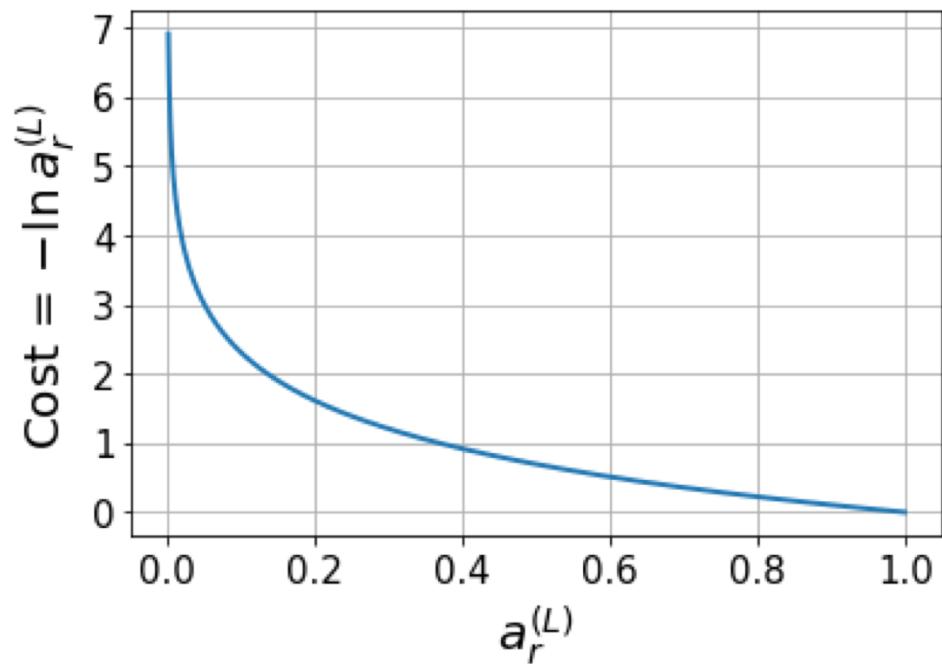
This makes beautiful sense as the error vector!

Example

$$a^{(L)} = \begin{bmatrix} 0.05 \\ 0.84 \\ 0.11 \\ 0 \end{bmatrix}$$

$$y^{(L)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\delta^{(L)} = \begin{bmatrix} 0.05 \\ -0.16 \\ 0.11 \\ 0 \end{bmatrix}$$



But remember: We are interested
in cost as a function of \mathbf{W}, \mathbf{b}

Quadratic Cost (MSE)

$$C = \sum_{i=1}^{N^{(L)}} (y_i - a_i)^2$$

Mainly used for regression, or cases where \mathbf{y} and \mathbf{a} are not probability distributions

Which cost function to use?

Use cross-entropy for classification tasks

It corresponds to Maximum Likelihood of a categorical distribution



Outline

- Activation functions
- Cost function
- **Optimizers**
- Regularization
- Parameter initialization
- Normalization
- Data handling
- Hyperparameter selection

Level sets

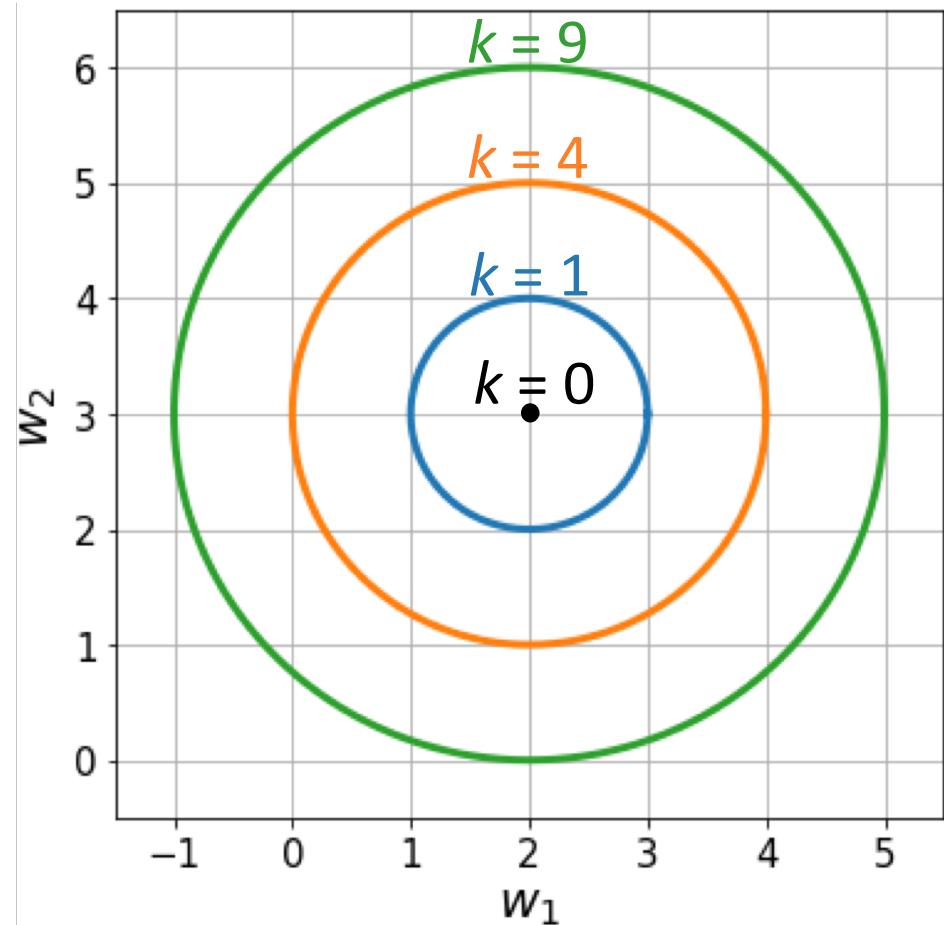
Given $f(\mathbf{w}) : \mathbb{R}^n \rightarrow \mathbb{R}$

A level set is a set of points in the domain at which function value is constant

$$f(\mathbf{w}) = k$$

Example

$$f(w_1, w_2) = (w_1 - 2)^2 + (w_2 - 3)^2$$

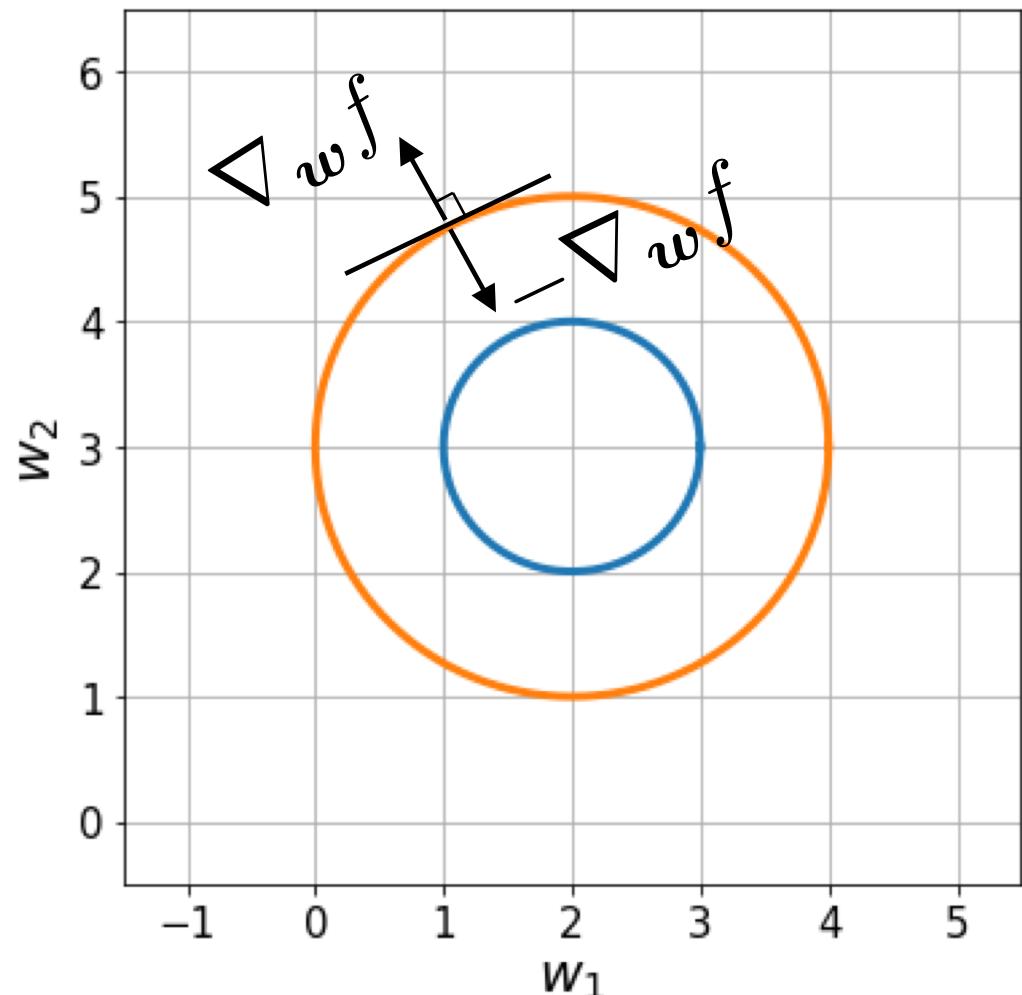


Level sets and gradient

Recall: Gradient at any point is the direction of maximum increase of the function at that point.

Gradient cannot have a component along the level set

The gradient (and its negative) are always perpendicular to the level set

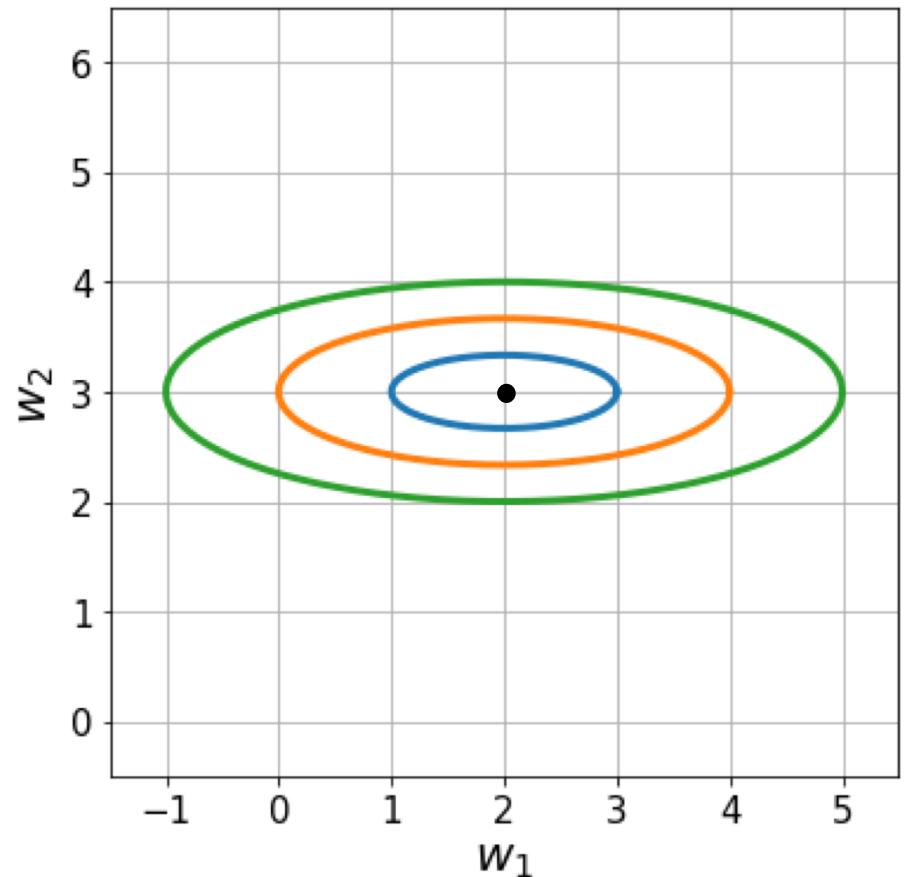


Conditioning

Informally, conditioning measures how non-circular the level sets are

Formally...

$$f(w_1, w_2) = (w_1 - 2)^2 + 9(w_2 - 3)^2$$



III-conditioned: Much more sensitive to w_2

Hessian: Matrix of double derivatives

Given $f(\mathbf{w}) : \mathbb{R}^n \rightarrow \mathbb{R}$

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \cdots & \frac{\partial^2 f}{\partial w_1 \partial w_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_m \partial w_1} & \cdots & \frac{\partial^2 f}{\partial w_m^2} \end{bmatrix}$$

For continuous functions:

$$\frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2 f}{\partial w_j \partial w_i}$$

Hessian is symmetric

Conditioning in terms of Hessian

Conditioning is measured as the condition number of the Hessian of the cost function

$$\kappa = \frac{\sigma_{\max}}{\sigma_{\min}} = \frac{|d_{\max}|}{|d_{\min}|}$$

I'm using d for eigenvalues

Examples:

$$C(w_1, w_2) = (w_1 - 2)^2 + (w_2 - 3)^2$$

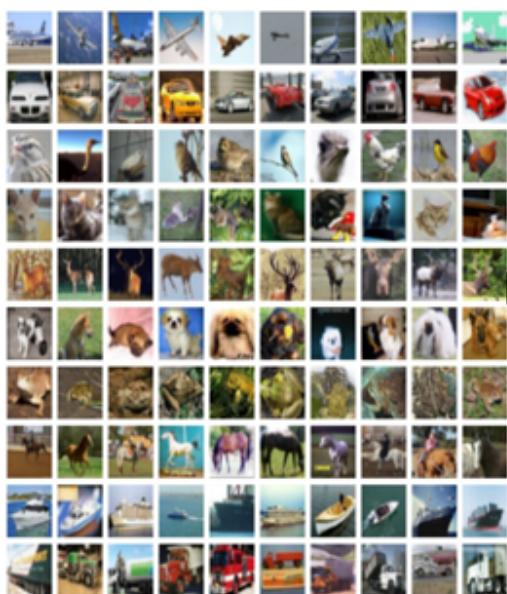
$\kappa = 1 \Rightarrow$ Well-conditioned

$$C(w_1, w_2) = (w_1 - 2)^2 + 9(w_2 - 3)^2$$

$\kappa = 9 \Rightarrow$ Ill-conditioned

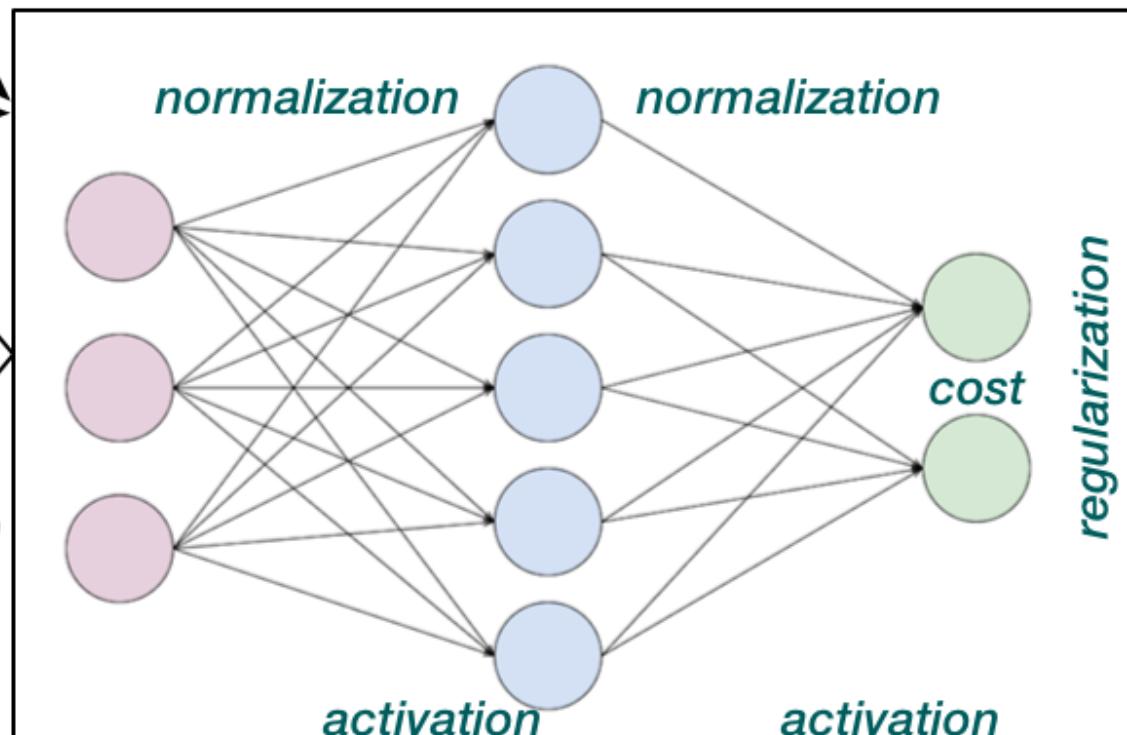
By the way, don't expect actual cost Hessians to be so simple...

The Big Picture



*data handling
normalization*

*parameter
initialization*



hyperparameter selection

What is an Optimizer?

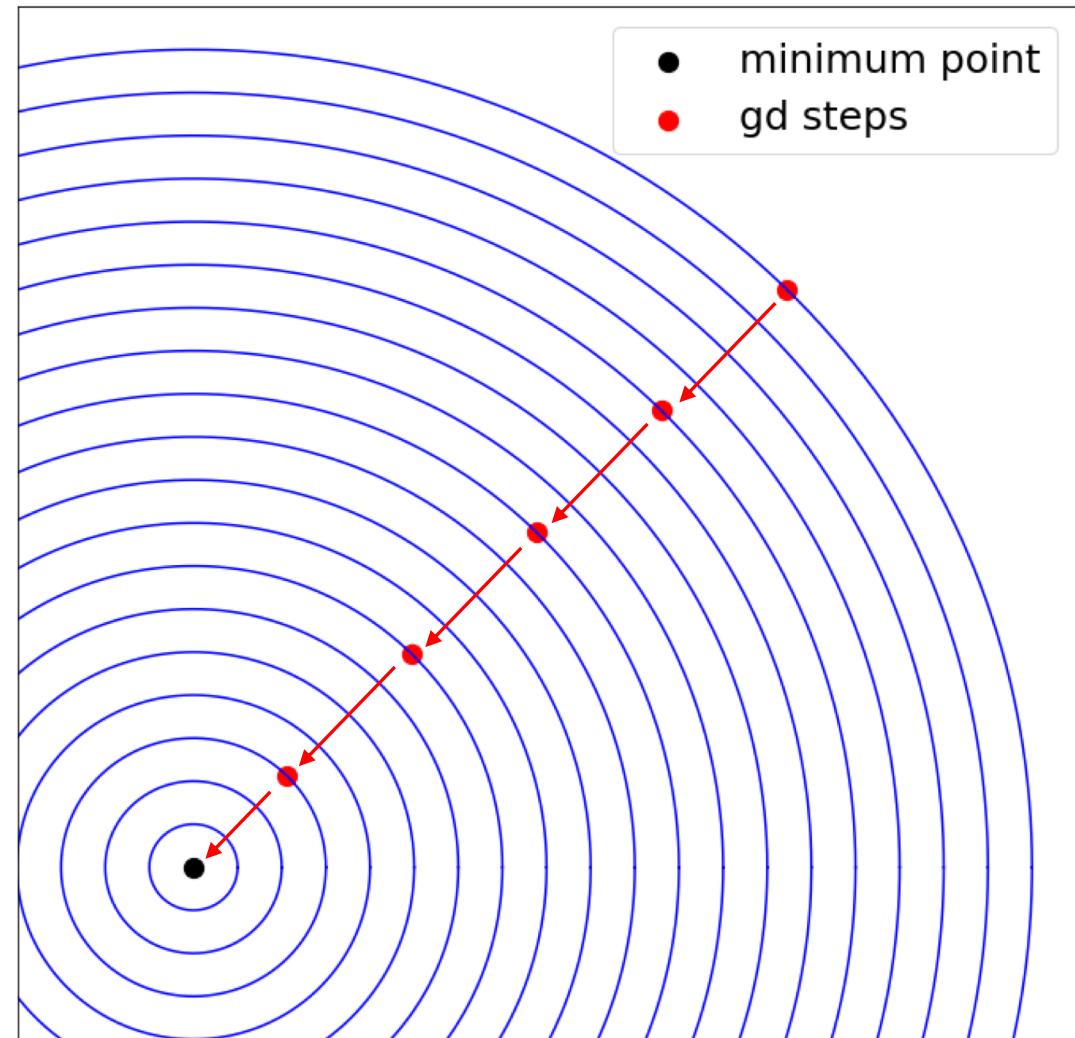
A strategy or algorithm to reduce the cost function and make the network learn

Eg: Gradient descent

$$p \leftarrow p - \eta \nabla_p C$$

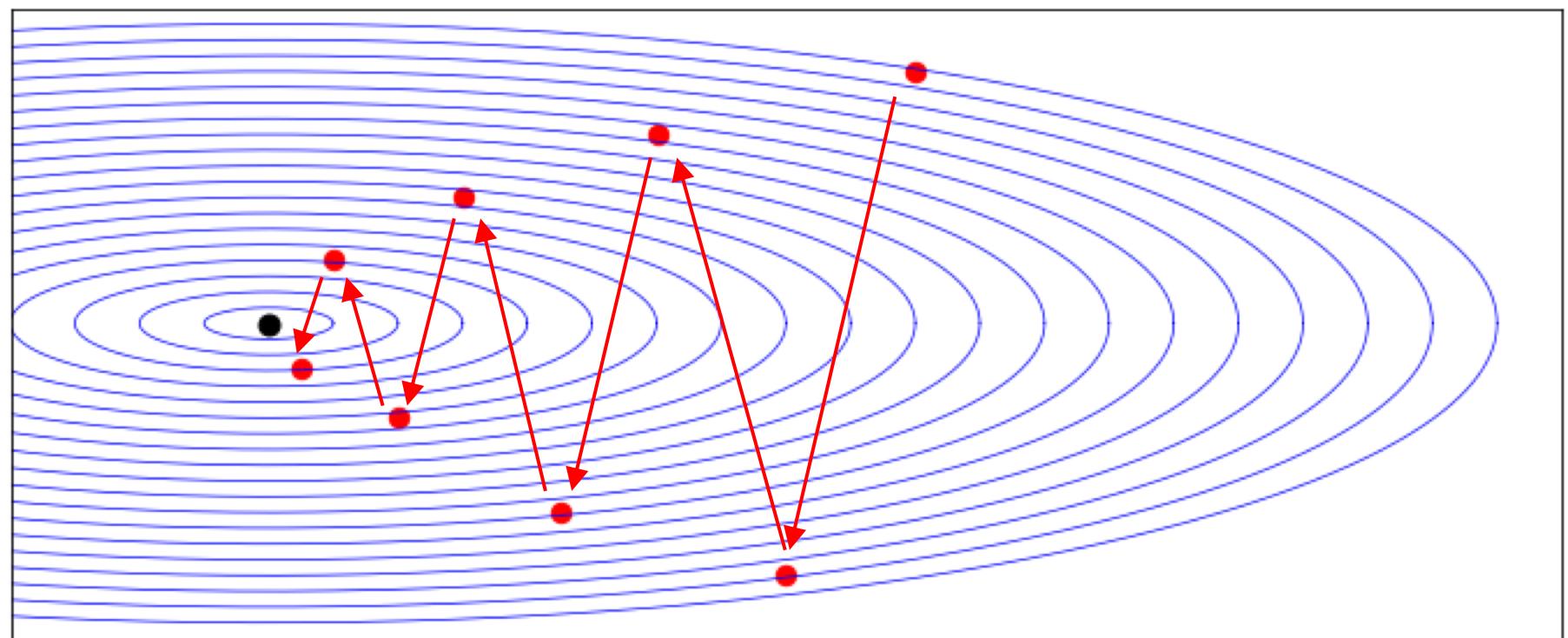
Gradient descent: Well-conditioned

*Well-conditioned level sets
lead to fast convergence*



Gradient descent: Ill-conditioned

*Ill-conditioned level sets
lead to slow convergence*



Momentum

Normal update:

$$p \leftarrow p - \eta \nabla_p C$$

Momentum update:

$$v \leftarrow \alpha v - \eta \nabla_p C$$

$$p \leftarrow p + v$$

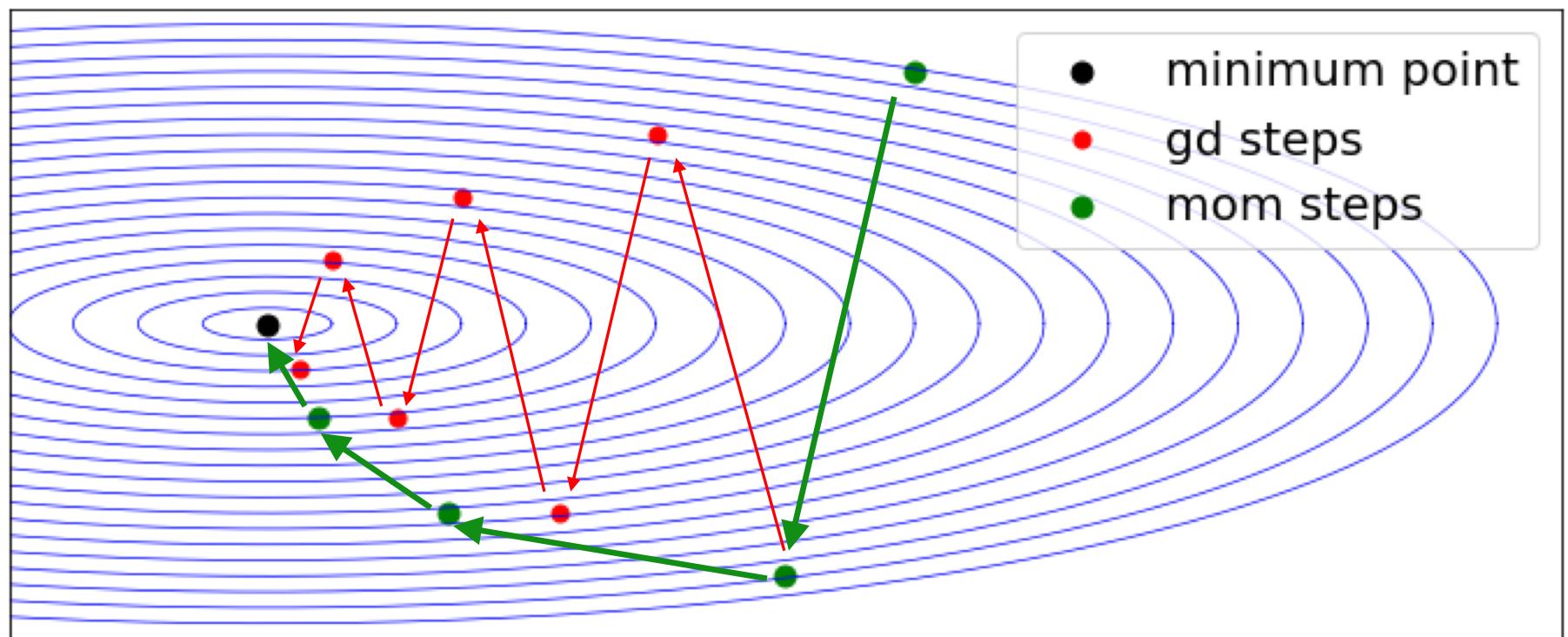
Equivalent to smoothing the update by a low pass filter

$$\alpha \in [0, 1]$$

How much of past history should be applied, typically ~ 0.9

Why Momentum?

Gradient descent with momentum
converges quickly even when ill-conditioned



Nesterov Momentum

Normal update:

$$p \leftarrow p - \eta \nabla_p C$$

Momentum update:

$$v \leftarrow \alpha v - \eta \nabla_p C$$

$$p \leftarrow p + v$$

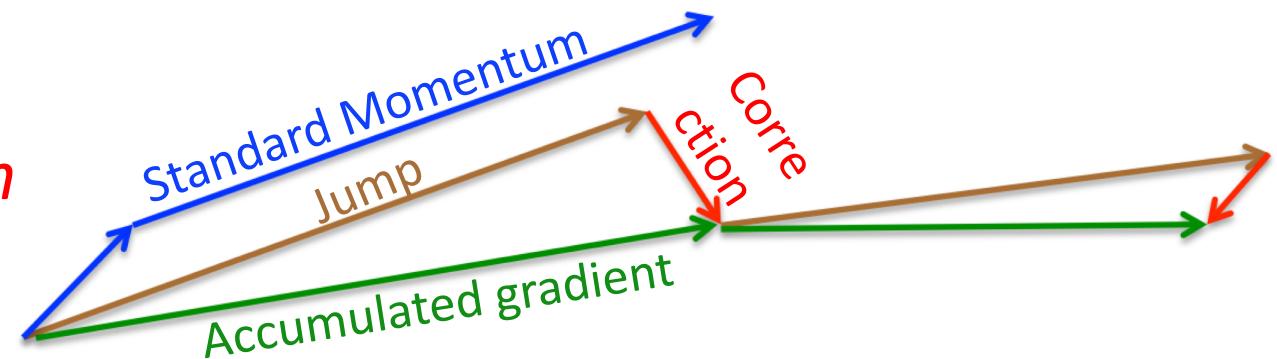
Nesterov Momentum update:

$$\tilde{p} \leftarrow p + \alpha v$$

$$v \leftarrow \alpha v - \eta \nabla_{\tilde{p}} C (\tilde{p})$$

$$p \leftarrow p + v$$

*Correction factor applied to momentum
May give faster convergence*



Adaptive Optimizers

- ✓ Learning rate should be reduced with time. We don't want to overshoot the minimum point and oscillate.
- ✓ Learning rate should be small for sensitive parameters, and vice-versa.
- ✓ Momentum factor should be increased with time to get smoother updates.

Adaptive optimizers:

- Adagrad
- Adadelta
- **RMSprop**
- **Adam**
- Adamax
- Nadam
- *Other 'Ada'-sounding words you can think of...*

RMSprop

Scale each gradient by an exponentially weighted moving average of its past history

$$v \leftarrow \rho v + (1 - \rho) (\nabla_p C)^2$$

$$p \leftarrow p - \frac{\eta}{\sqrt{v} + \epsilon} \nabla_p C$$

Default $\rho = 0.9$

ϵ is just a small number to prevent division by 0. Can be machine epsilon

Sensitive parameters get low η

Adam

RMSprop with momentum and bias correction

At time step (t+1):

$$v_1 \leftarrow \rho_1 v_1 + (1 - \rho_1) \nabla_p C$$

Momentum

$$v_2 \leftarrow \rho_2 v_2 + (1 - \rho_2) (\nabla_p C)^2$$

Exponentially weighted
moving average of past history

$$\tilde{v}_1 = \frac{v_1}{1 - \rho_1^t}$$
$$\tilde{v}_2 = \frac{v_2}{1 - \rho_2^t}$$

Bias corrections to make

$$\mathbb{E}\{\tilde{v}_1\} = \mathbb{E}\{\nabla_p C\}$$

$$\mathbb{E}\{\tilde{v}_2\} = \mathbb{E}\{(\nabla_p C)^2\}$$

$$p \leftarrow p - \frac{\eta}{\sqrt{\tilde{v}_2} + \epsilon} \tilde{v}_1$$

Defaults:

$$\eta=0.001, \rho_1 = 0.9, \rho_2 = 0.999$$

Which optimizer to use?



`opt = SGD(eta=0.01, mom=0.95, decay=1e-5, nesterov=False)`

`opt = Adam()`

Machine learning vs Simple Optimization

	Simple Optimization	Machine Learning
Goal	Minimize $f(x)$	Minimize $C(w)$ on test data
Typical problem size	A few variables	A few million variables
Approach	Gradient descent 2^{nd} order methods	Gradient descent on training data (2^{nd} order methods not feasible)
Stopping criterion	$x^* = \operatorname{argmin} f(x)$?? No access to test data

Machine learning is about generalizing well on test data

Outline

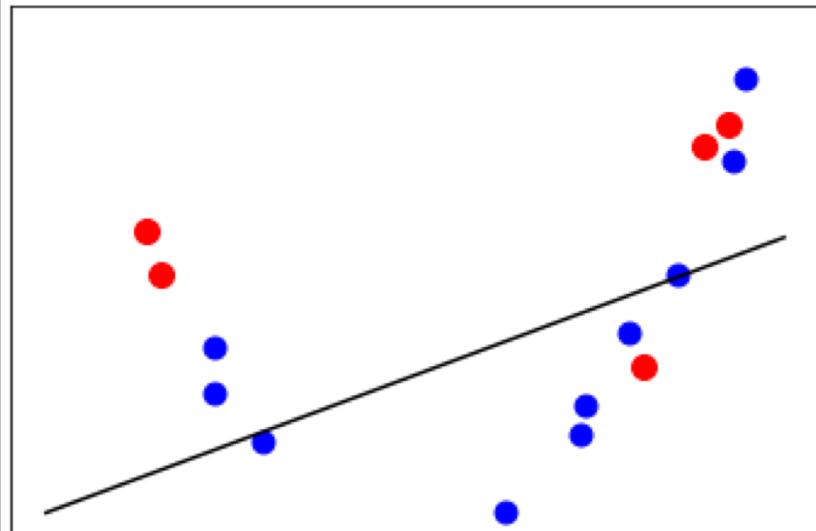
- Activation functions
- Cost function
- Optimizers
- **Regularization**
- Parameter initialization
- Normalization
- Data handling
- Hyperparameter selection

Regularization

Regularization is any modification intended to reduce generalization (test) cost, but not training cost

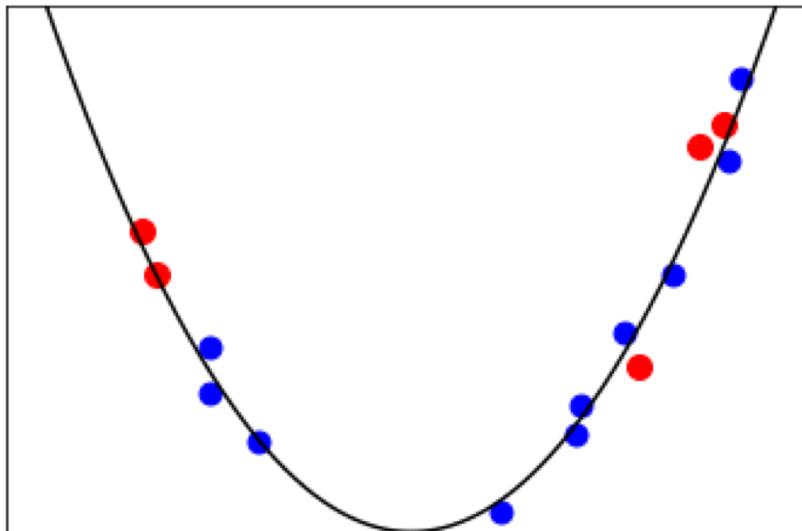
The problem of overfitting

Underfitting: Order = 1



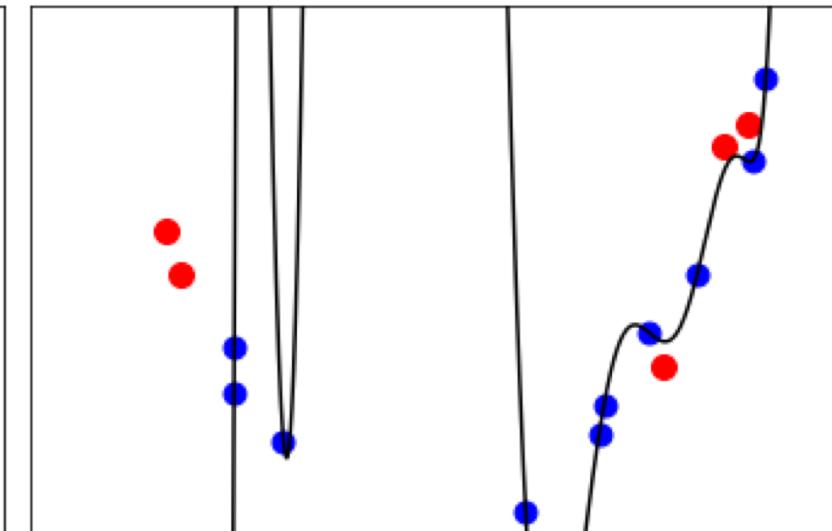
$$C_{\text{train}} = 0.385$$
$$C_{\text{test}} = 0.49$$

Perfect: Order = 2



$$C_{\text{train}} = 0.012$$
$$C_{\text{test}} = 0.019$$

Overfitting: Order = 9



$$C_{\text{train}} = 0$$
$$C_{\text{test}} = 145475$$



Bias-Variance Tradeoff

Say the (unknown) optimum value of some parameter is w^* , but the value we get from training is w

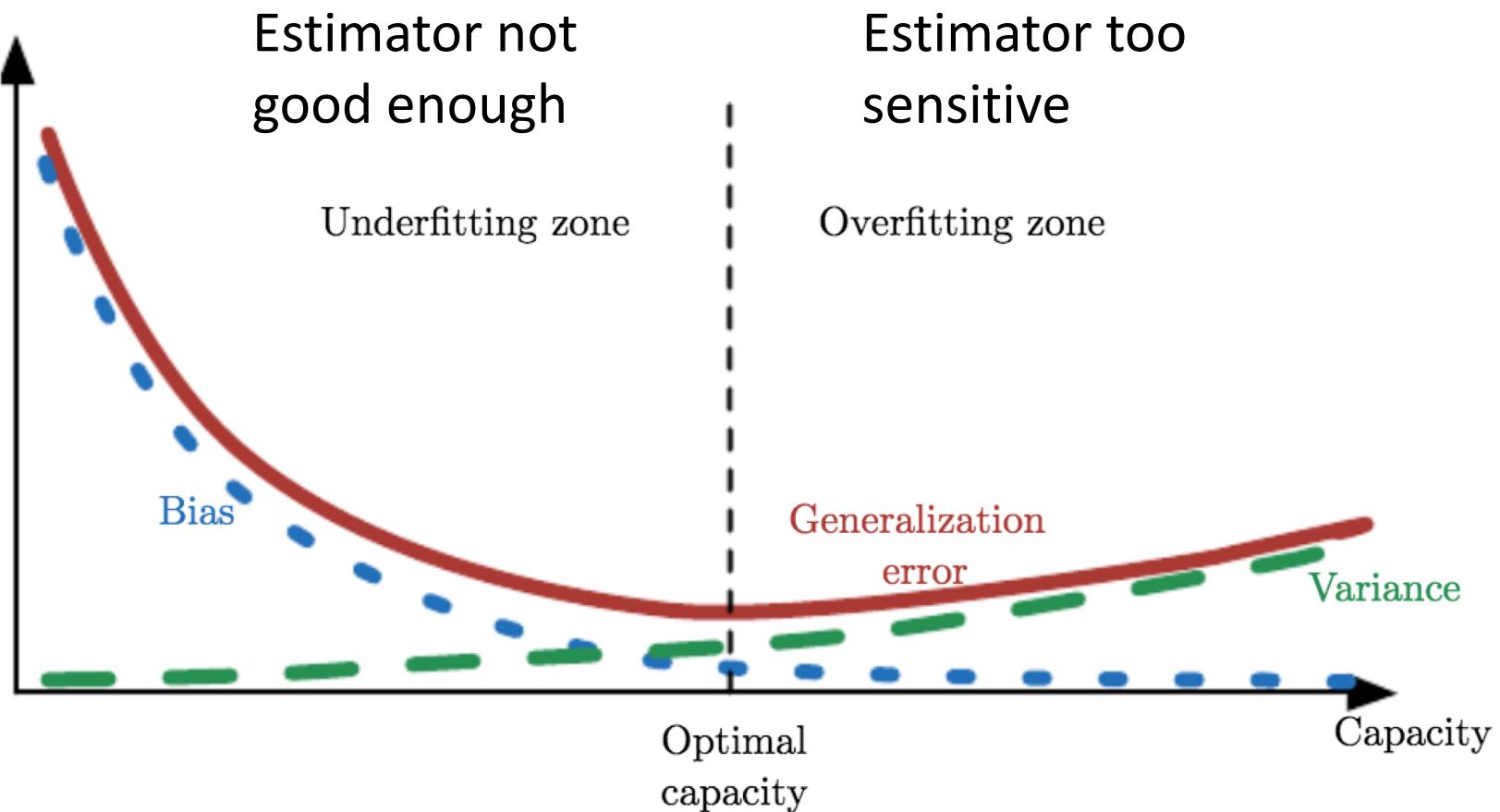
$$\begin{aligned}\text{MSE} &= \mathbb{E} \left\{ (w - w^*)^2 \right\} \\ &= \mathbb{E} \{ w^2 \} - 2w^* \mathbb{E}\{w\} + w^{*2} \\ &= \mathbb{E} \{ w^2 \} - \mathbb{E}\{w\}^2 + \mathbb{E}\{w\}^2 - 2w^* \mathbb{E}\{w\} + w^{*2} \\ &= (\mathbb{E} \{ w^2 \} - \mathbb{E}\{w\}^2) + (\mathbb{E}(w) - w^*)^2 \\ &= \text{Var}(w) + \text{Bias}(w)^2\end{aligned}$$

Reduce Bias: Make expected value of estimate close to true value

Reduce Variance: Estimate bounds should be tight

For a given MSE, bias trades off with variance

Bias-Variance Tradeoff



Minibatches

$$p \leftarrow p - \eta \langle \nabla_p C \rangle_{\text{traindata}}$$

*Estimate true gradient using
data average of sample
gradients of training data*

$$p \leftarrow p - \frac{\eta}{N_{\text{train}}} \sum_{i=1}^{N_{\text{train}}} (\nabla_p C)^{[i]}$$

Batch gradient descent
Best possible approximation to true gradient

$$p \leftarrow p - \frac{\eta}{M} \sum_{i=1}^M (\nabla_p C)^{[i]}$$

Minibatch gradient descent
Noisy approximation to true gradient

A Note on Terminology

Batch size : Same as minibatch size M

Batch GD : Use all samples

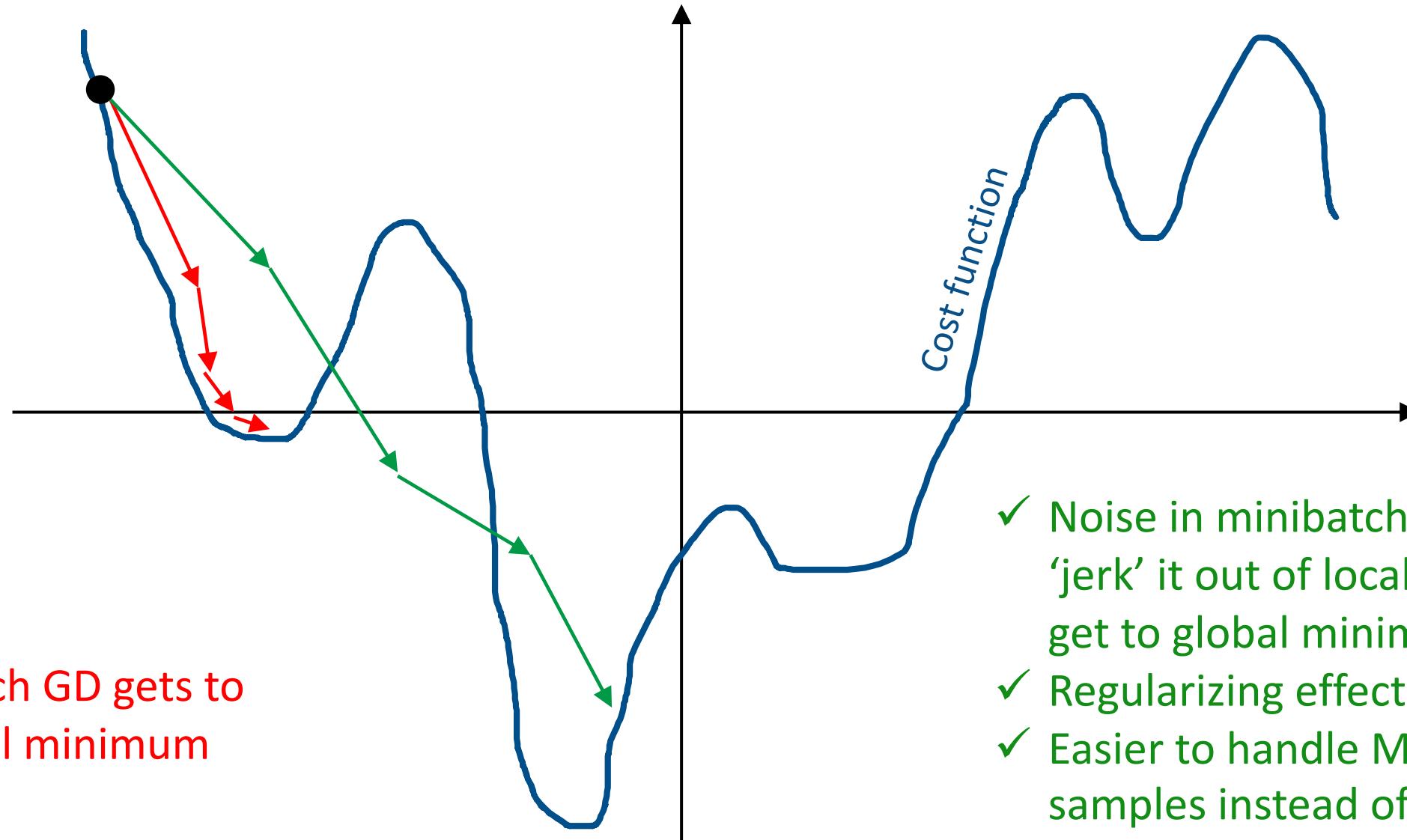
Minibatch GD : $1 < M < N_{\text{train}}$

Online learning : $M = 1$

Stochastic GD : $1 \leq M < N_{\text{train}}$

Basically the same as minibatch GD.

Optimizers - Why minibatches?



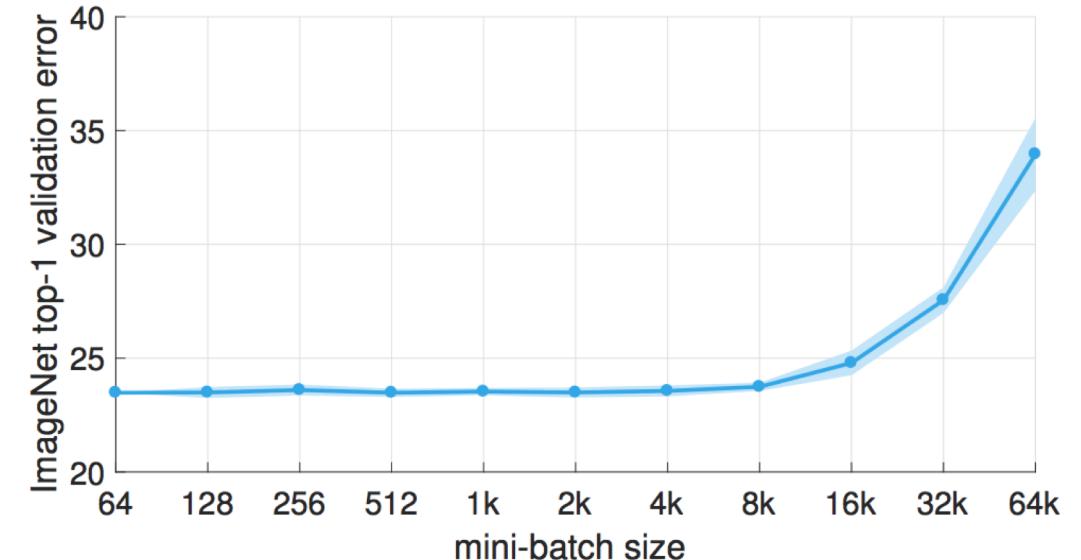
Choosing a batch size

Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, Kaiming He. "Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour". arXiv:1706.02677

$M \leq 8000$ is fine

But...

Dominic Masters, Carlo Luschi. "Revisiting Small Batch Training for Deep Neural Networks". arXiv:1804.07612



Recommend $M = 16, 32, 64$

Both use the same dataset: Imagenet, with $N_{train} > 1$ million

So what batch size do I choose?

No consensus



- ✓ Depends on the number of workers/threads in CPU/GPU
- ✓ Powers of 2 work well
- ✓ Keras default is 32 : Good for small datasets ($N_{\text{train}} < 10000$)
- ✓ For bigger datasets, I would recommend trying 128, 256, even 1024
- ✓ Your computer might slow down with $M > 1000$

Concept of Penalty Function

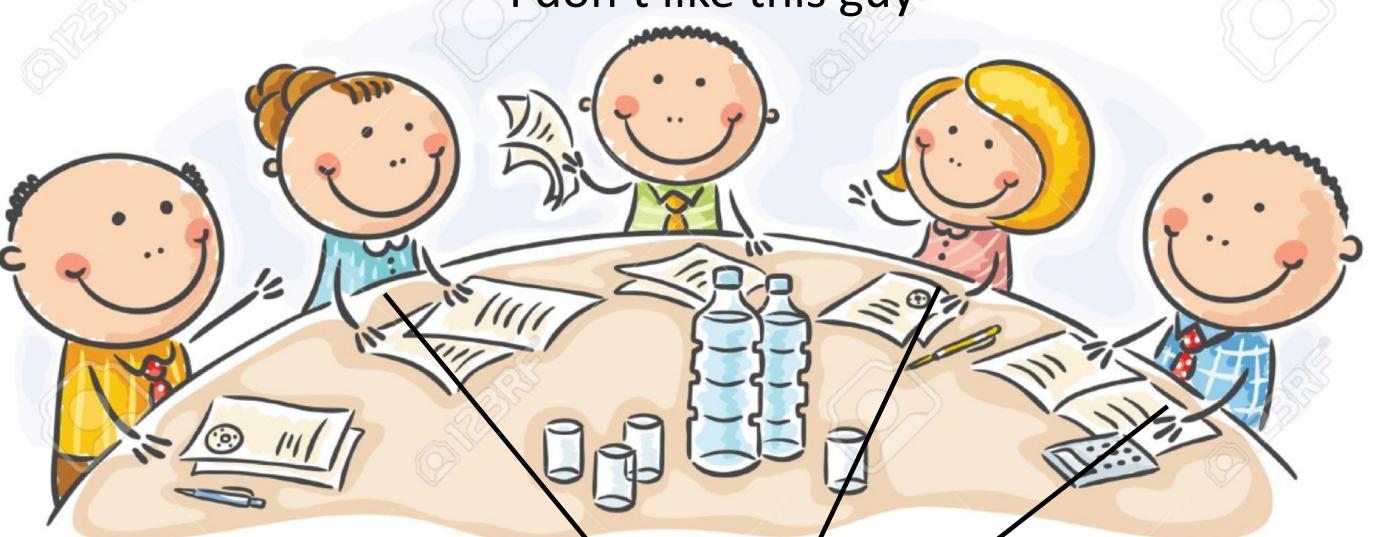
Penalty measures dislike

Eg: Penalty = Yellow dress

Big penalty
I absolutely
hate this guy

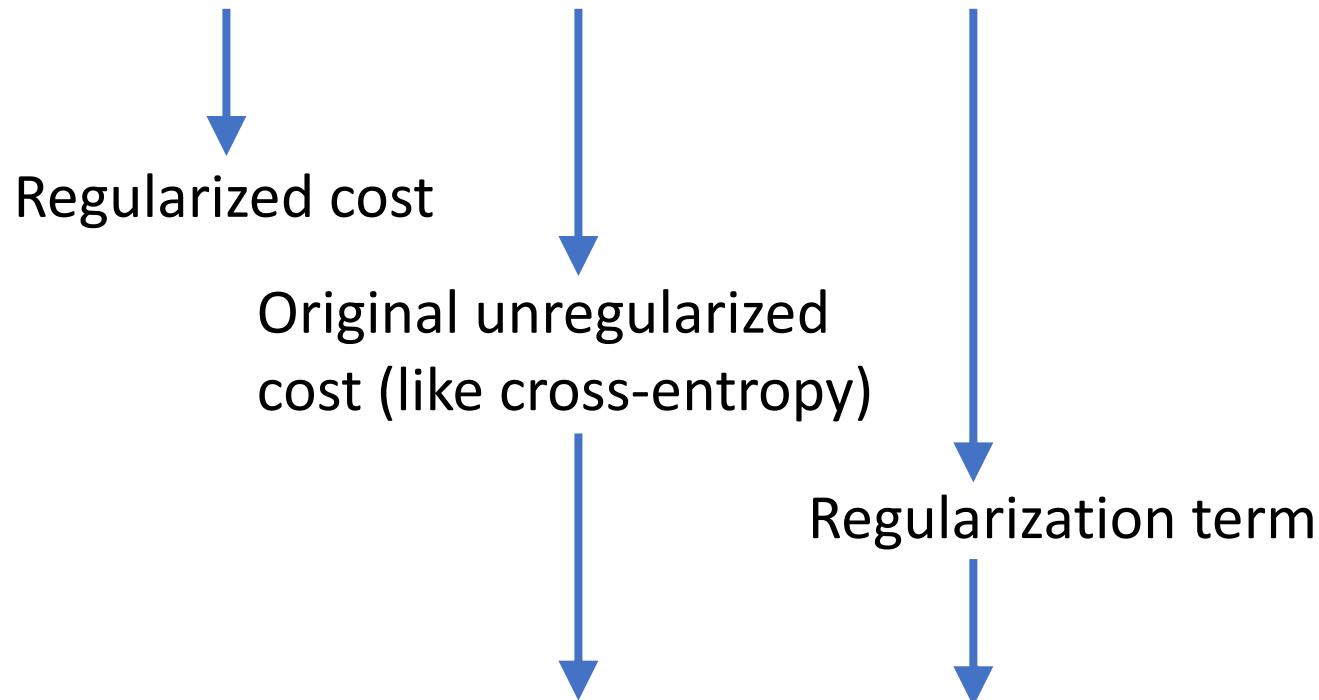
Small penalty
I don't like this guy

No penalty
These people are fine



L2 weight penalty (Ridge)

$$C(\mathbf{w}) = C_0(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$



Update rule: $\mathbf{w} \leftarrow \mathbf{w} - \eta (\nabla_{\mathbf{w}} C_0 + 2\lambda \mathbf{w})$

Regularization hyperparameter effect

$$\lambda = \frac{\text{Importance of small weights}}{\text{Importance of reducing training cost}}$$

$\lambda = 0$ Reduce training cost only

$$\boldsymbol{w}^* = \arg \min C_0(\boldsymbol{w})$$

Overfitting

$\lambda = \infty$ Reduce weight norm only

$$\boldsymbol{w}^* = \mathbf{0}$$

Underfitting

Typically $\lambda \in [10^{-5}, 10^{-3}]$

Effect on Conditioning

Consider a 2nd order Taylor approximation of the unregularized cost function C_0 around its optimum point \mathbf{w}^* , i.e. $\nabla C_0(\mathbf{w}^*) = 0$

$$C_0(\mathbf{w}) = C_0(\mathbf{w}^*) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^T \mathbf{H} (\mathbf{w} - \mathbf{w}^*)$$

$$\nabla_{\mathbf{w}} C_0(\mathbf{w}) = \mathbf{H} (\mathbf{w} - \mathbf{w}^*)$$

$$\nabla_{\mathbf{w}} C(\mathbf{w}) = \mathbf{H} (\mathbf{w} - \mathbf{w}^*) + 2\lambda \mathbf{w}$$


Adding the gradient of the regularization term to get gradient of the regularized cost function C

Effect on Conditioning

Set the gradient to 0 to find optimum $\tilde{\mathbf{w}}$ of the regularized cost

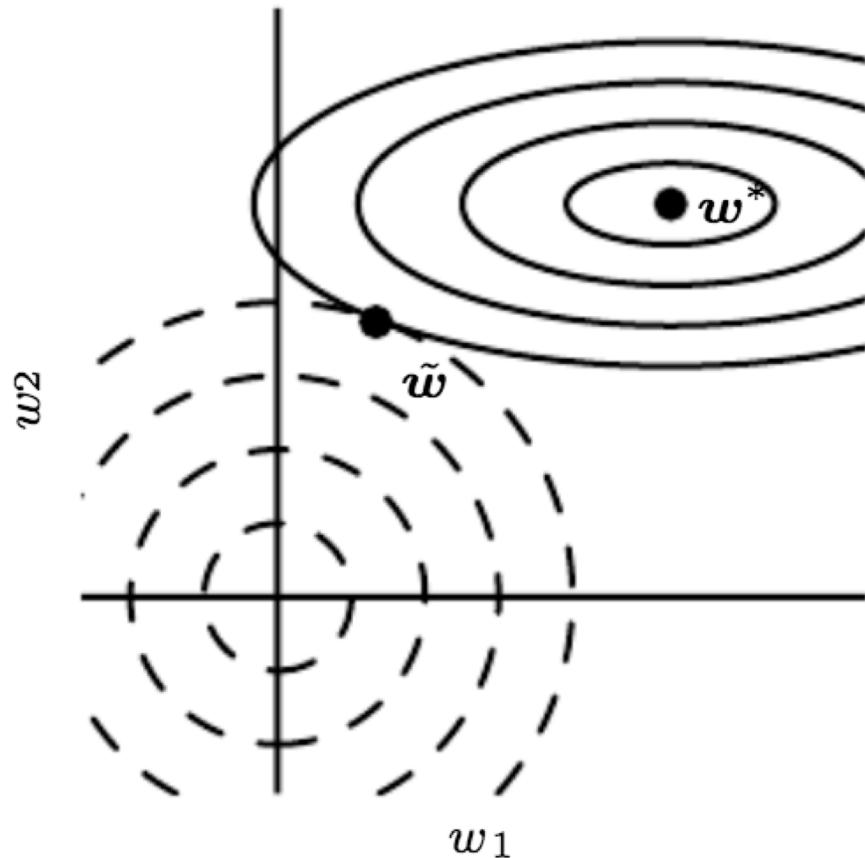
$$\nabla_{\mathbf{w}} C(\mathbf{w}) = \mathbf{0}$$

$$\begin{aligned}\tilde{\mathbf{w}} &= (\mathbf{H} + 2\lambda\mathbf{I})^{-1} \mathbf{H} \mathbf{w}^* \\ &= (\mathbf{Q} \mathbf{D} \mathbf{Q}^T + 2\lambda\mathbf{I})^{-1} \mathbf{Q} \mathbf{D} \mathbf{Q}^T \mathbf{w}^* \\ &= [\mathbf{Q}(\mathbf{D} + 2\lambda\mathbf{I})\mathbf{Q}^T]^{-1} \mathbf{Q} \mathbf{D} \mathbf{Q}^T \mathbf{w}^* \\ &= \mathbf{Q} \left[(\mathbf{D} + 2\lambda\mathbf{I})^{-1} \mathbf{D} \right] \mathbf{Q}^T \mathbf{w}^*\end{aligned}$$

Eigendecomposition of \mathbf{H}
 \mathbf{D} is the diagonal matrix
of eigenvalues

Any eigenvalue d_i of \mathbf{H} is replaced by $d_i / (d_i + 2\lambda)$ => Improves conditioning

Effect on Weights



— Level sets of C_0
- - - Level sets of norm penalty

Less important weights like w_1 are driven towards 0 by regularization

L1 weight penalty (LASSO)

$$C(\mathbf{w}) = C_0(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

Regularized cost

Original unregularized
cost (like cross-entropy)

Regularization term

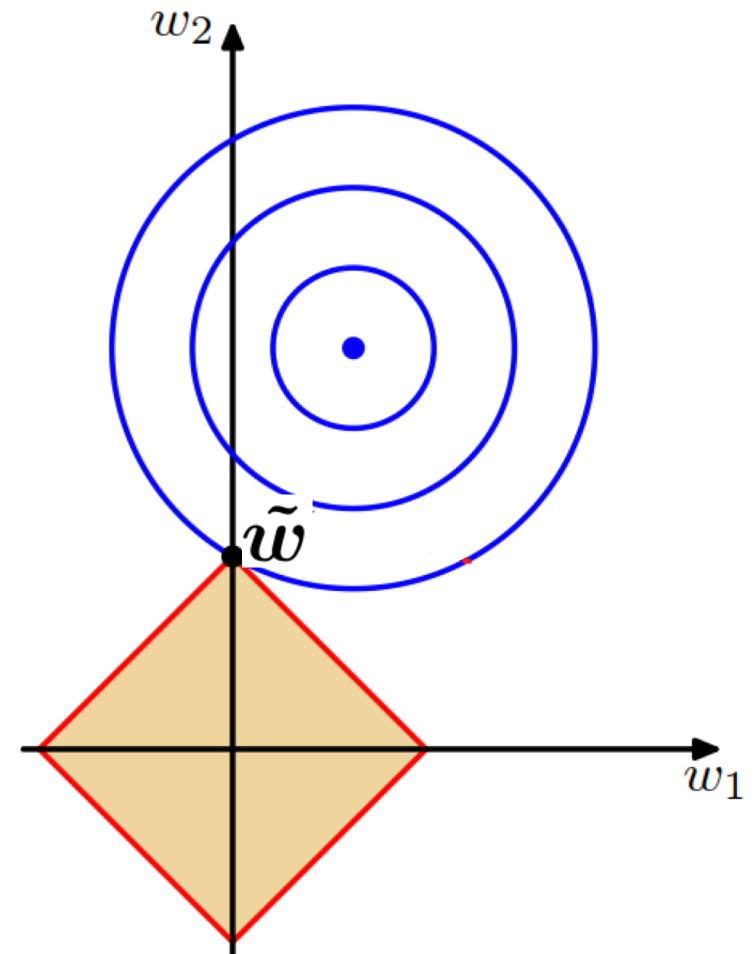
Update rule: $\mathbf{w} \leftarrow \mathbf{w} - \eta (\nabla_{\mathbf{w}} C_0 + \lambda \text{sign}(\mathbf{w}))$

Sparsity

L1 promotes sparsity => Drives weights to 0

Eg: Some weight $w = 0.1$

	Function value	Dislike	Update value	Action
L2	$w^2 = 0.01$	Small	$\eta\lambda w = 0.1(\eta\lambda)$	Leave it
L1	$ w = 0.1$	Large	$\eta\lambda * \text{sign}(w) = \eta\lambda$	Make it smaller



Other weight penalties

L0 : Count number of nonzero weights.

This is a great way to get sparsity, but not differentiable and hard to analyze

Elastic Net - Combine L1 and L2 :

Reduce some weights, zero out some

Also see Group lasso:

W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li, "Learning structured sparsity in deep neural networks," in Proc. Advances in Neural Information Processing Systems 29 (NIPS), 2016, pp. 2074–2082.

Weight penalty as MAP

$$\begin{aligned}\tilde{\mathbf{w}} &= \arg \max p(\mathbf{w}|\text{Data}) \\&= \arg \max p(\text{Data}|\mathbf{w})p(\mathbf{w}) \\&= \arg \max [\log p(\text{Data}|\mathbf{w}) + \log p(\mathbf{w})] \\&= \arg \min [-\log p(\text{Data}|\mathbf{w})] \quad + \quad \arg \min [-\log p(\mathbf{w})]\end{aligned}$$



Original cost C_0



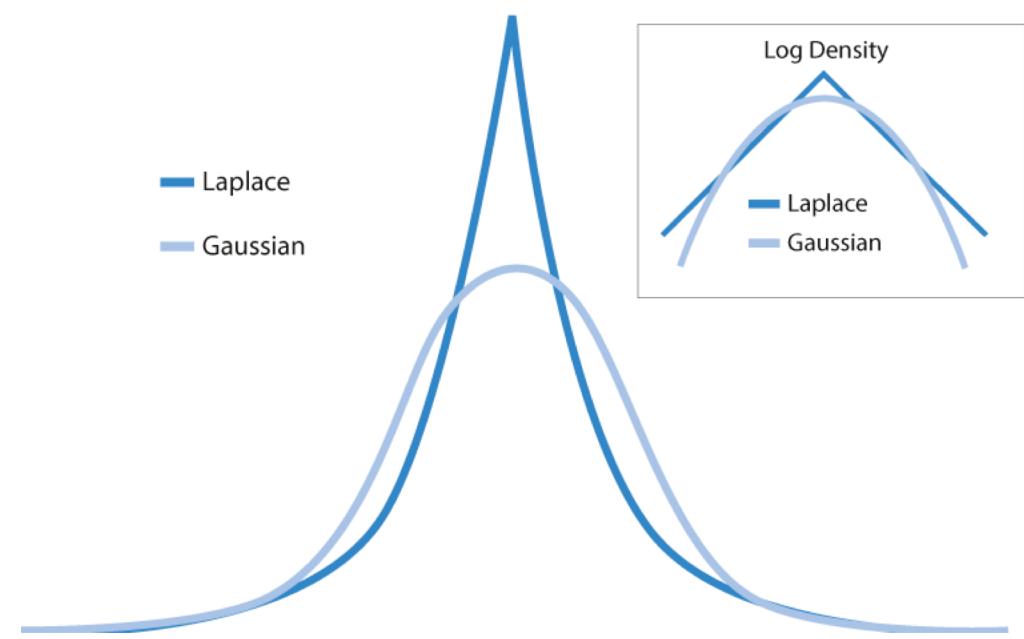
Weight penalty

Weight penalty is our prior belief about the weights

Weight penalty as MAP

$$-\log \left[p(\mathbf{w}) = \mathcal{N} \left(\mathbf{w}, \mathbf{0}, \frac{1}{\lambda} \mathbf{I} \right) \right] = \lambda \|\mathbf{w}\|_2^2$$

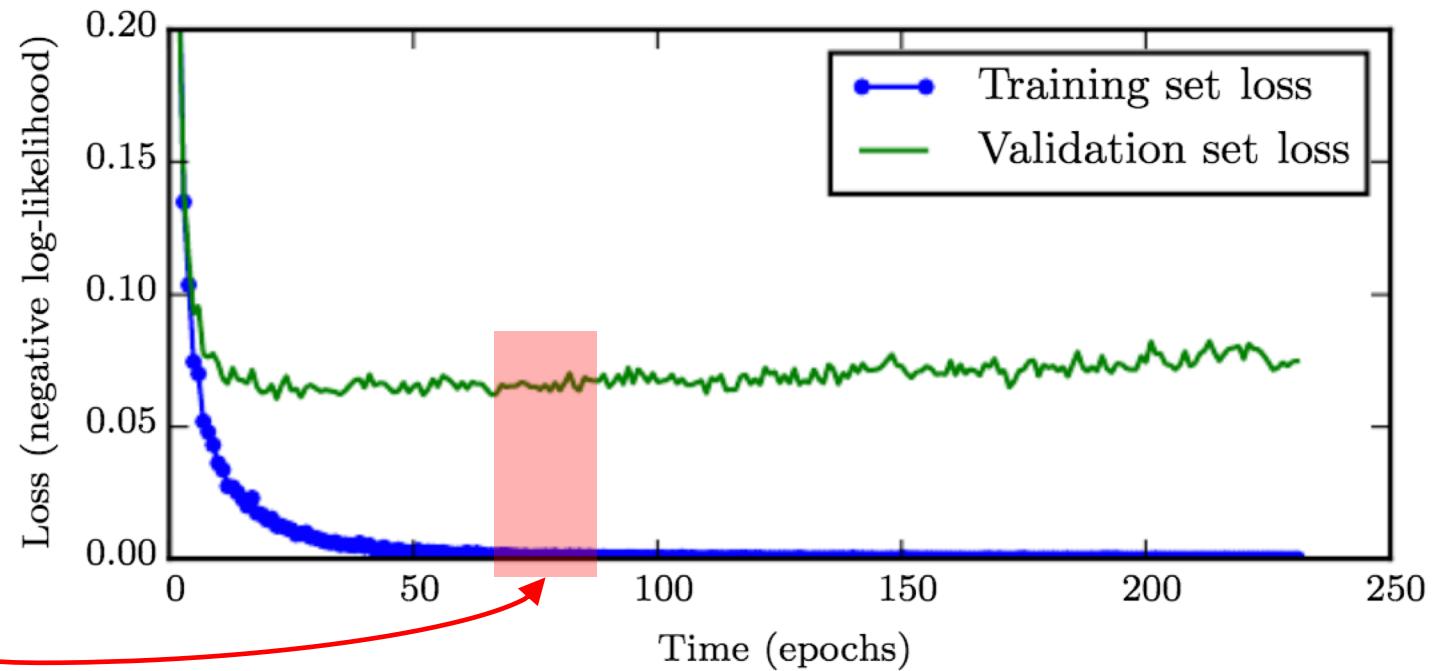
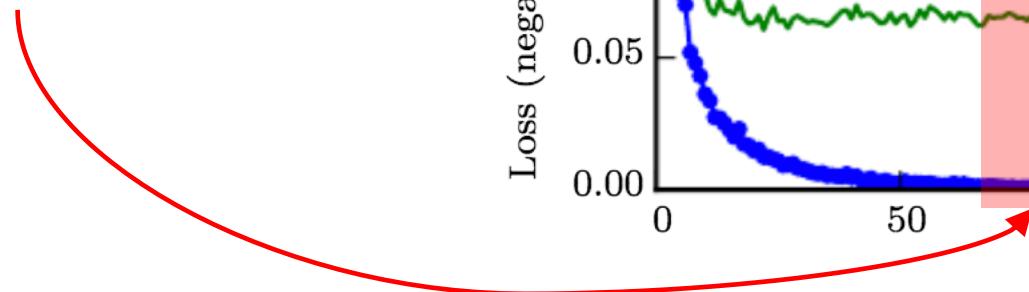
$$-\log \left[p(w_i) = \text{Laplace} \left(w_i, 0, \frac{1}{\lambda} \right) \right] = \lambda \|w\|_1$$



Early stopping

Validation metrics are used as a proxy for test metrics

Stop training if validation metrics don't get better

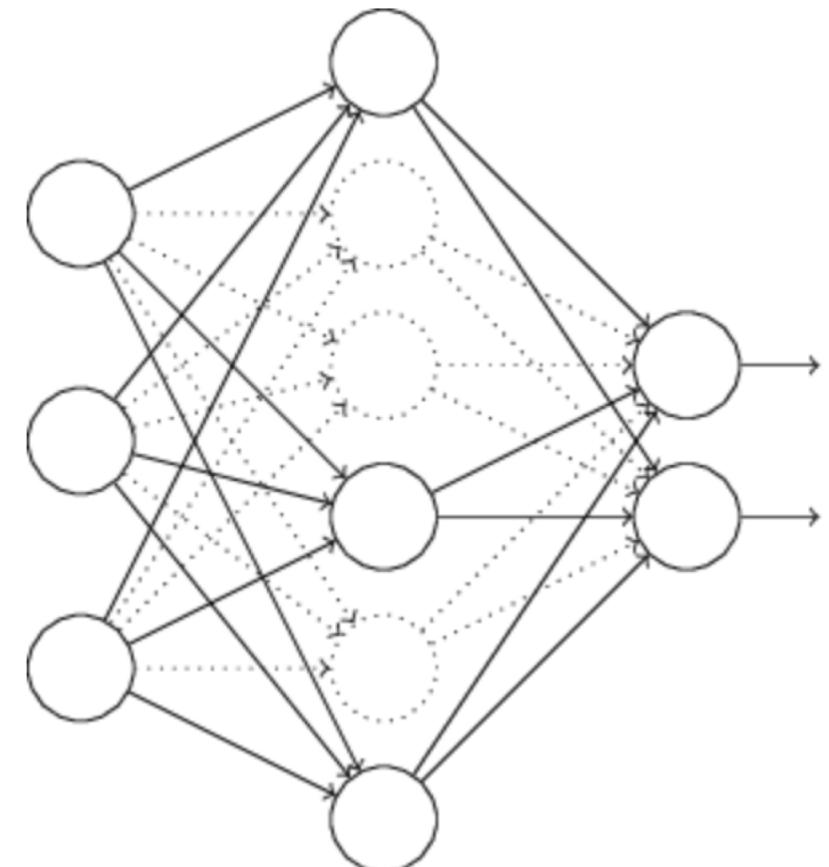


Dropout

For every minibatch during training, delete a random selection of input and hidden nodes

Keep probability p_k : Fraction of units to keep (i.e. not drop)

But why dropout??



Ensemble Methods

Let's train an ensemble of k similar networks on the same problem and average their test results
(Random differences: Initialization, dataset shuffling)

Error of a single network: ϵ_i

Error of the ensemble: $\frac{1}{k} \sum_i \epsilon_i$

$$\begin{aligned}\text{MSE}_{\text{ens}} &= \mathbb{E} \left\{ \left(\frac{1}{k} \sum_i \epsilon_i \right)^2 \right\} \\ &= \frac{1}{k^2} \mathbb{E} \left\{ \sum_i \left(\epsilon_i^2 + \sum_{j \neq i} \epsilon_i \epsilon_j \right) \right\} \\ &= \frac{1}{k} \mathbb{E} \{ \epsilon_i^2 \} + \frac{k-1}{k} \mathbb{E} \{ \epsilon_i \epsilon_j \}\end{aligned}$$

Averaging the results of multiple networks reduces errors

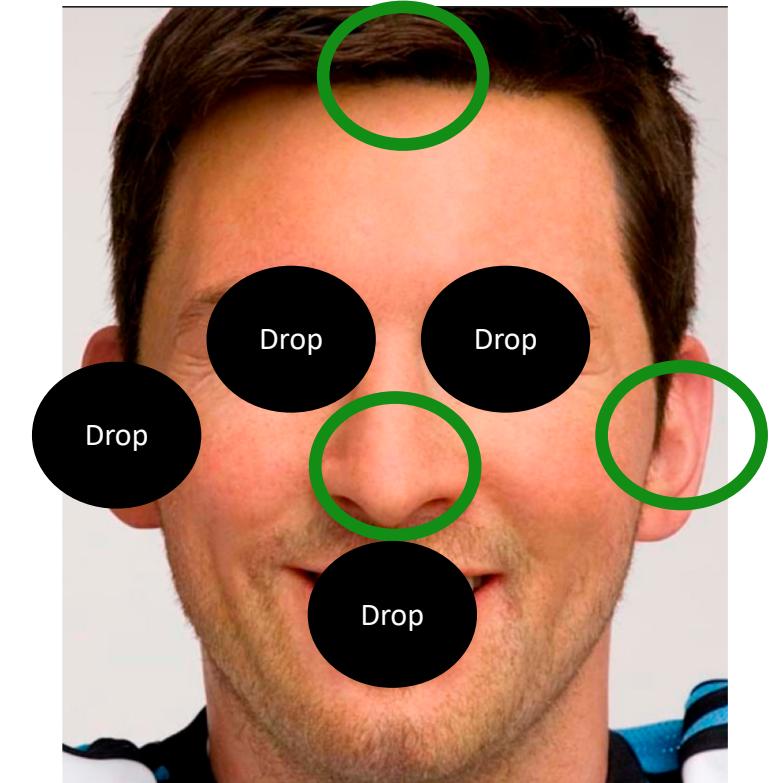
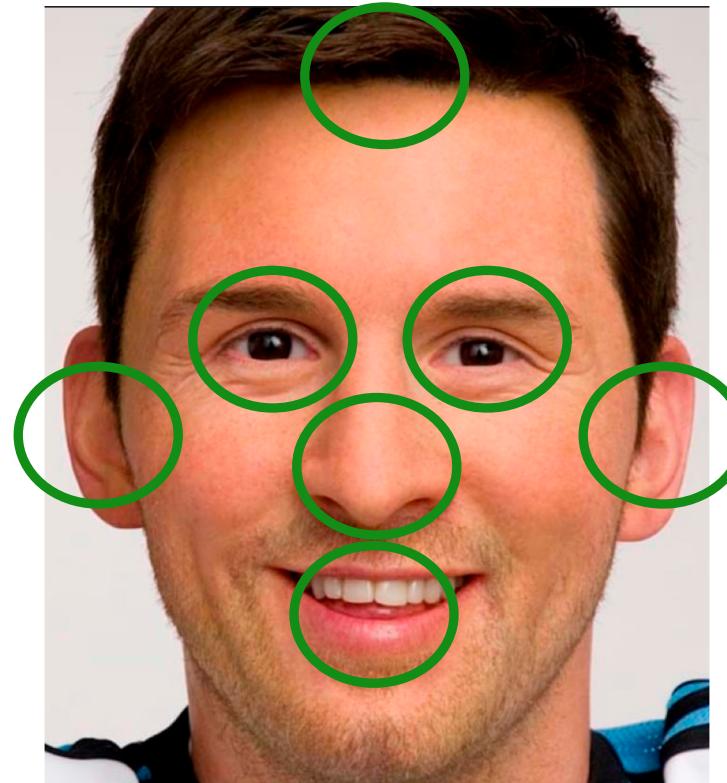
$$\begin{aligned}\because \mathbb{E} \{ \epsilon_i \epsilon_j \} &\leq \mathbb{E} \{ \epsilon_i^2 \} \\ \Rightarrow \text{MSE}_{\text{ens}} &\leq \text{MSE}\end{aligned}$$

But this process is expensive!

Return to dropout

Dropout is an ensemble method
using the same shared weights
=> Computationally cheaper

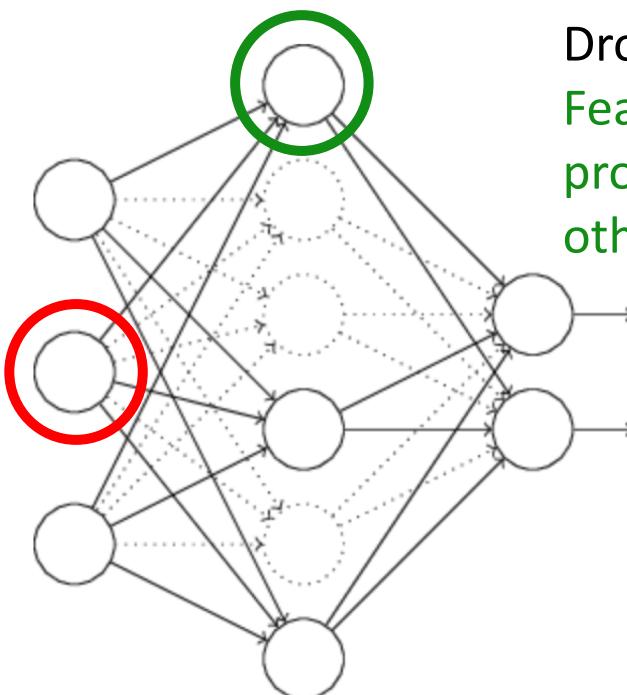
*Improves robustness, since
individual nodes have to work
without depending on other
deleted nodes => Regularization*



How much dropout?

Original paper recommendations:

Input $p_k = 0.8$



Hidden $p_k = 0.5$

Drop hidden:
Features still
propagate through
other hidden nodes

0.5 gives maximum regularization effect:
P. Baldi, P. J. Sadowski, "Understanding Dropout", Proc. NIPS 2013

But you should try other values as well!

Output $p_k = 1$

Drop output: Cannot classify!

Drop input: Feature gone forever!
Some people use input $p_k = 1$

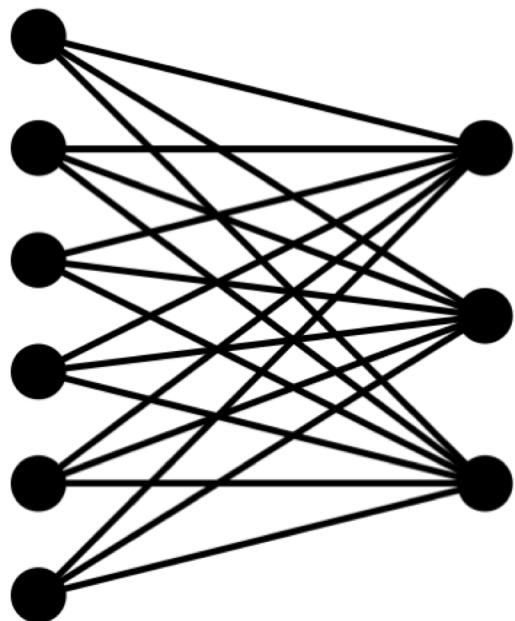
Conv layer $p_k = 0.7$

N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov,
"Dropout: A simple way to prevent neural networks from overfitting,"
Journal of Machine Learning Research, vol. 15, pp. 1929–1958, 2014

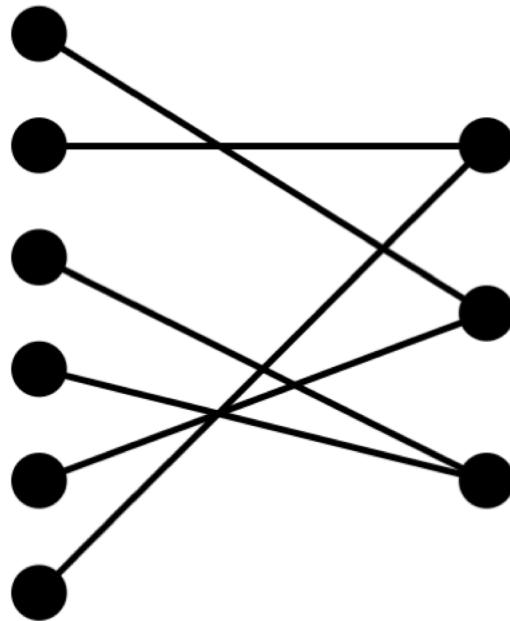
More on Dropout

- All nodes and weights are present in inference. Multiply weights by p_k during inference to compensate
 - This assumes linearity, but still works!
- Dropout is like multiplicative 0-1 noise for each node
- Dropout is best for big networks
 - Increase capacity, increase generalization power, but prevent overfitting

My Research: Pre-defined sparsity



Fully-connected network



Pre-defined sparse network

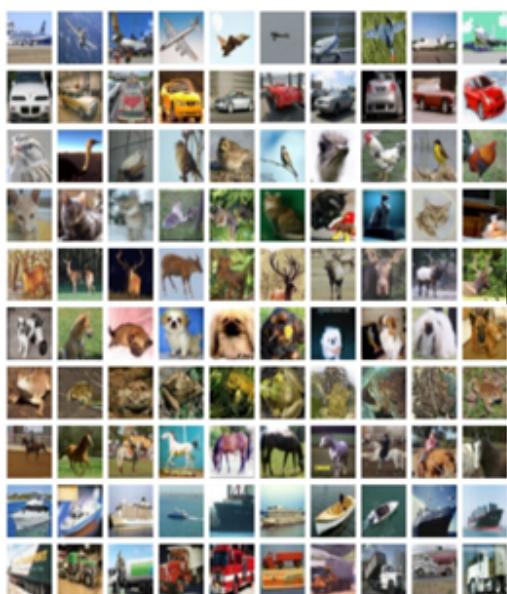
*Remove weights from
the very beginning
Train and test on a low
complexity network*

How to design a good
connection pattern?

Outline

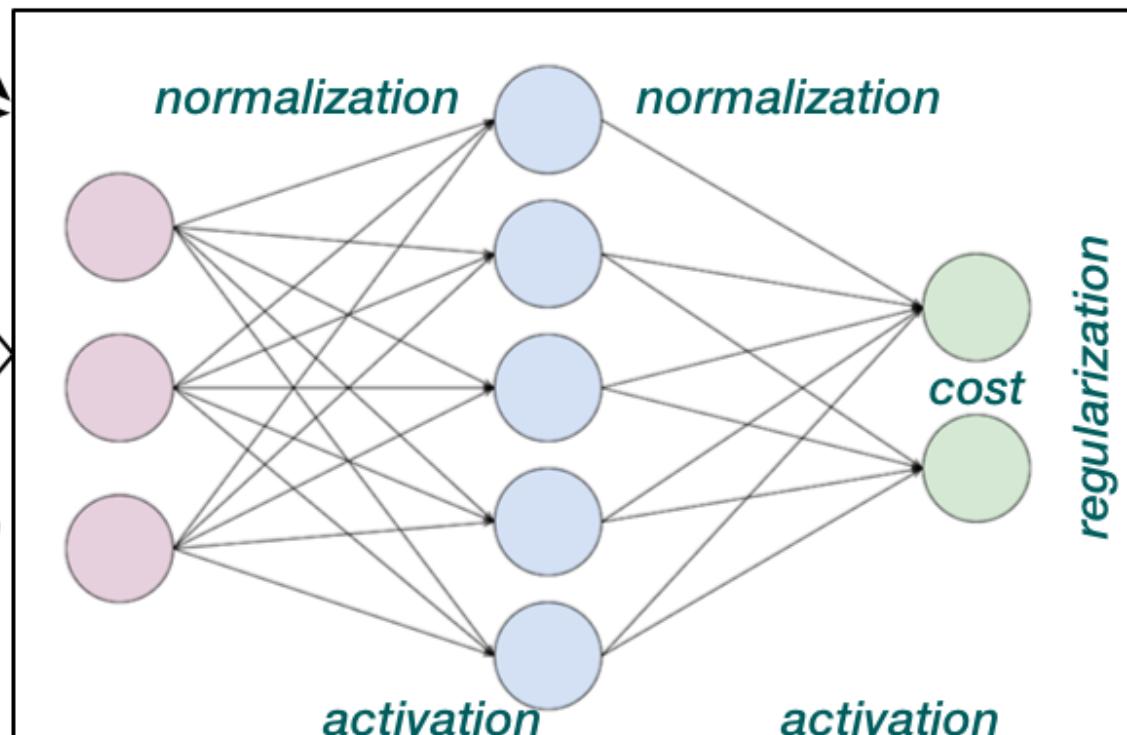
- Activation functions
- Cost function
- Optimizers
- Regularization
- **Parameter initialization**
- Normalization
- Data handling
- Hyperparameter selection

The Big Picture



*data handling
normalization*

*parameter
initialization*



hyperparameter selection

Parameter Initialization

$$p \leftarrow p - \eta \nabla_p C$$

What is the initial value of p ?

How about.... Same initialization, like all 0s

BAD

We want each weight to be useful on its own, not mirror other weights

Think of a company, do 2 people have exactly the same job?

Instead, initialize weights randomly

Glorot (Xavier) Normal Initialization

Imagine a linear function: $y = w_1x_1 + w_2x_2 + \cdots + w_Nx_N$

Say all w, x are IID: $\text{Var}(y) = N\text{Var}(w)\text{Var}(x)$

$$\text{if } \text{Var}(w) = \frac{1}{N}$$

$$\text{then } \text{Var}(y) = \text{Var}(x)$$

For forward prop : $\text{Var}(W^{(l)}) = 1/N^{(l-1)}$

For backprop : $\text{Var}(W^{(l)}) = 1/N^{(l)}$

Compromise :

$$W^{(l)} \sim \mathcal{N}\left(0, \frac{2}{N^{(l)} + N^{(l-1)}}\right)$$

Despite a lot of assumptions, it works well!!

Xavier Glorot, Yoshua Bengio. Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, PMLR 9:249-256, 2010.

Why...

... should variances be equal?

... should we take harmonic mean as compromise?

The original paper doesn't explicitly explain these

I think of it as similar to vanishing and exploding gradient issues.

From a forward-propagation point of view, to keep information flowing we would like that

$$\forall(i, i'), \text{Var}[z^i] = \text{Var}[z^{i'}]. \quad (8)$$

From a back-propagation point of view we would similarly like to have

$$\forall(i, i'), \text{Var}\left[\frac{\partial \text{Cost}}{\partial s^i}\right] = \text{Var}\left[\frac{\partial \text{Cost}}{\partial s^{i'}}\right]. \quad (9)$$

These two conditions transform to:

$$\forall i, n_i \text{Var}[W^i] = 1 \quad (10)$$

$$\forall i, n_{i+1} \text{Var}[W^i] = 1 \quad (11)$$

As a compromise between these two constraints, we might want to have

$$\forall i, \text{Var}[W^i] = \frac{2}{n_i + n_{i+1}} \quad (12)$$

Glorot (Xavier) Uniform Initialization

Let $W^{(l)} \sim U(-a, a)$

$$\Rightarrow \text{Var}(W^{(l)}) = \frac{a^2}{3}$$

$$\Rightarrow \frac{a^2}{3} = \frac{2}{N^{(l)} + N^{(l-1)}}$$

$$W^{(l)} \sim U\left(-\sqrt{\frac{6}{N^{(l)} + N^{(l-1)}}}, \sqrt{\frac{6}{N^{(l)} + N^{(l-1)}}}\right)$$

He Initialization

Glorot normal assumes linear functions

Actual activations like ReLU are non-linear

So...

$$W^{(l)} \sim \mathcal{N} \left(0, \frac{2}{N^{(l-1)}} \right)$$

$$W^{(l)} \sim U \left(-\sqrt{\frac{6}{N^{(l-1)}}}, \sqrt{\frac{6}{N^{(l-1)}}} \right)$$

Bias Initialization and Choosing an Initializer

He initialization works slightly better than Glorot

Nothing to choose between Normal and Uniform

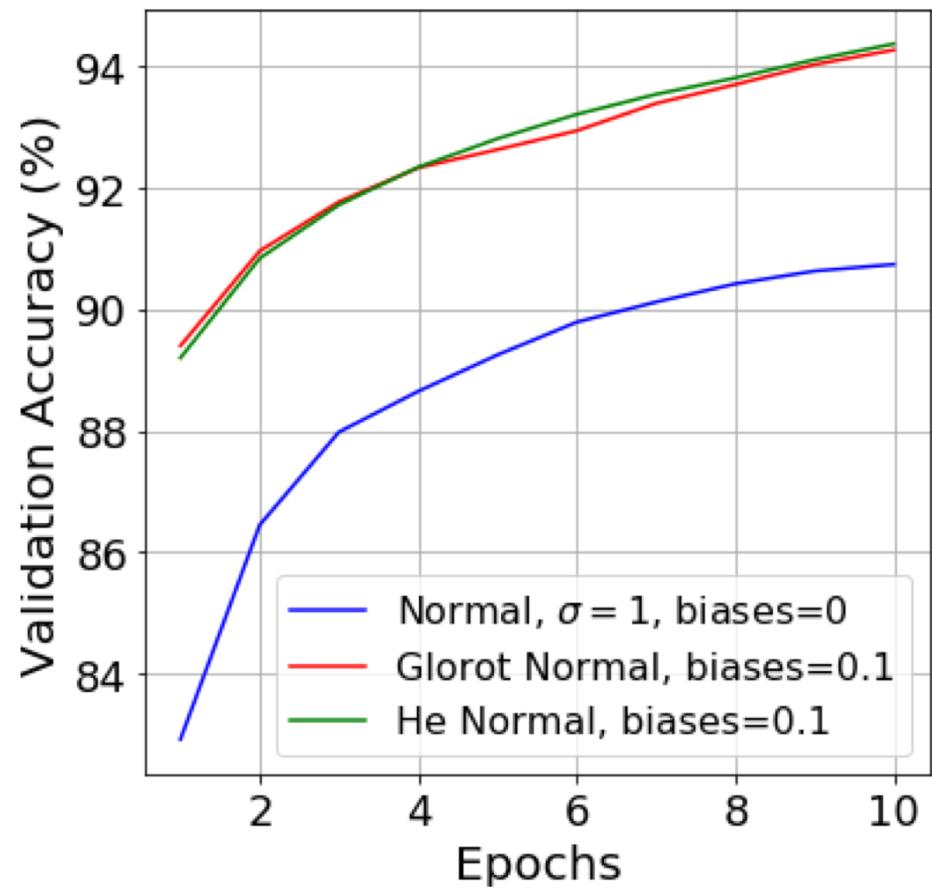
Biases not as important as weights

All zeros work OK

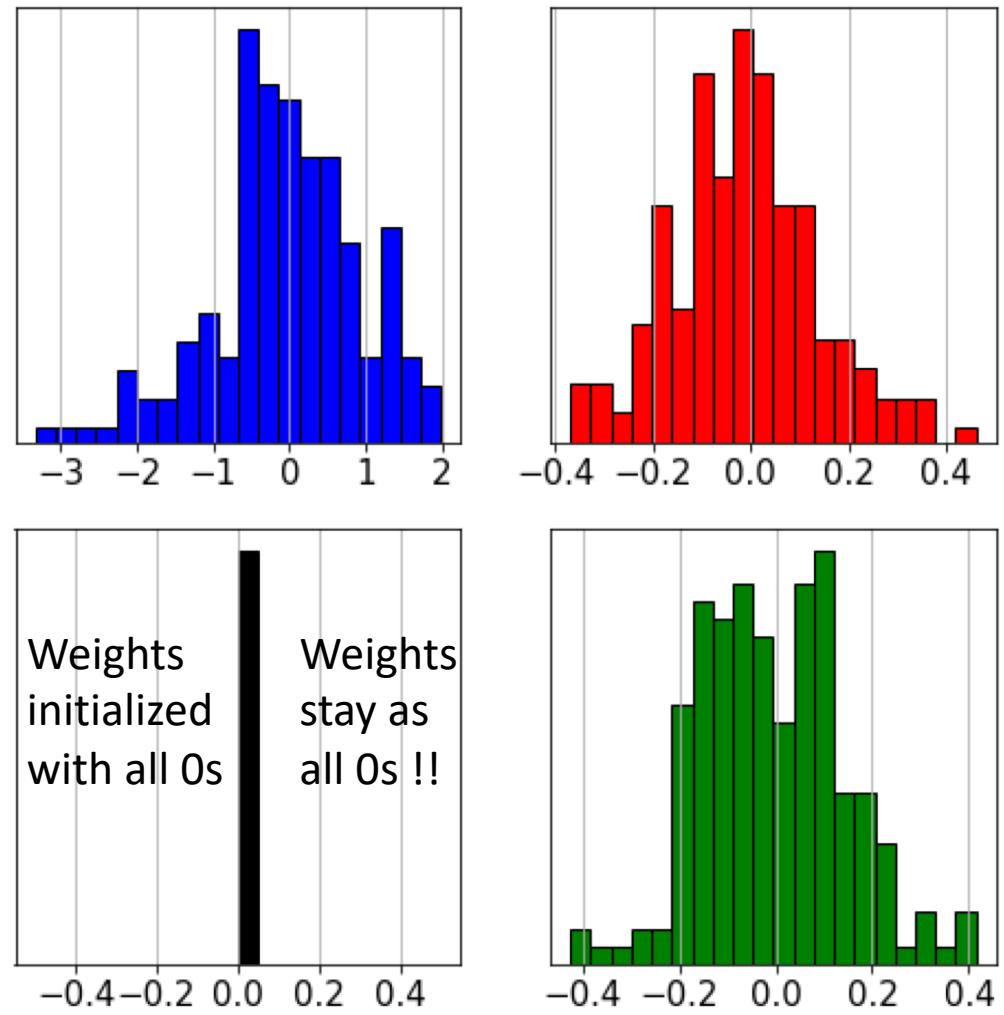
For ReLU, give a small positive value like 0.1
to biases to prevent neurons dying



Comparison of Initializers



MNIST [784, 200, 10]
Regularization: None



Histograms of a few weights in 2nd junction after training for 10 epochs

Outline

- Activation functions
- Cost function
- Optimizers
- Regularization
- Parameter initialization
- **Normalization**
- Data handling
- Hyperparameter selection

Normalization

Imagine trying to predict how good a footballer is...

Feature	Units	Range
Height	Meters	1.5 to 2
Weight	Kilograms	50 to 100
Shot speed	Kmph	120 to 180
Shot curve	Degrees	0 to 10
Age	Years	20 to 35
Minutes played	Minutes	5,000 to 20,000
<i>Fake diving?</i>	--	Yes / No

Different features have very different scales



Input normalization (Preprocessing)

Say there are n total training samples, each with f features

Eg for MNIST: $n=50000, f=784$

$$\mathbf{X}_{\text{train}} = \begin{bmatrix} x_{11} & \cdots & x_{1f} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nf} \end{bmatrix}$$

*Compute stats for each feature
across all training samples*

μ_i : Mean

σ_i : Standard deviation

M_i : Maximum value

m_i : Minimum value

for all i from 1 to f

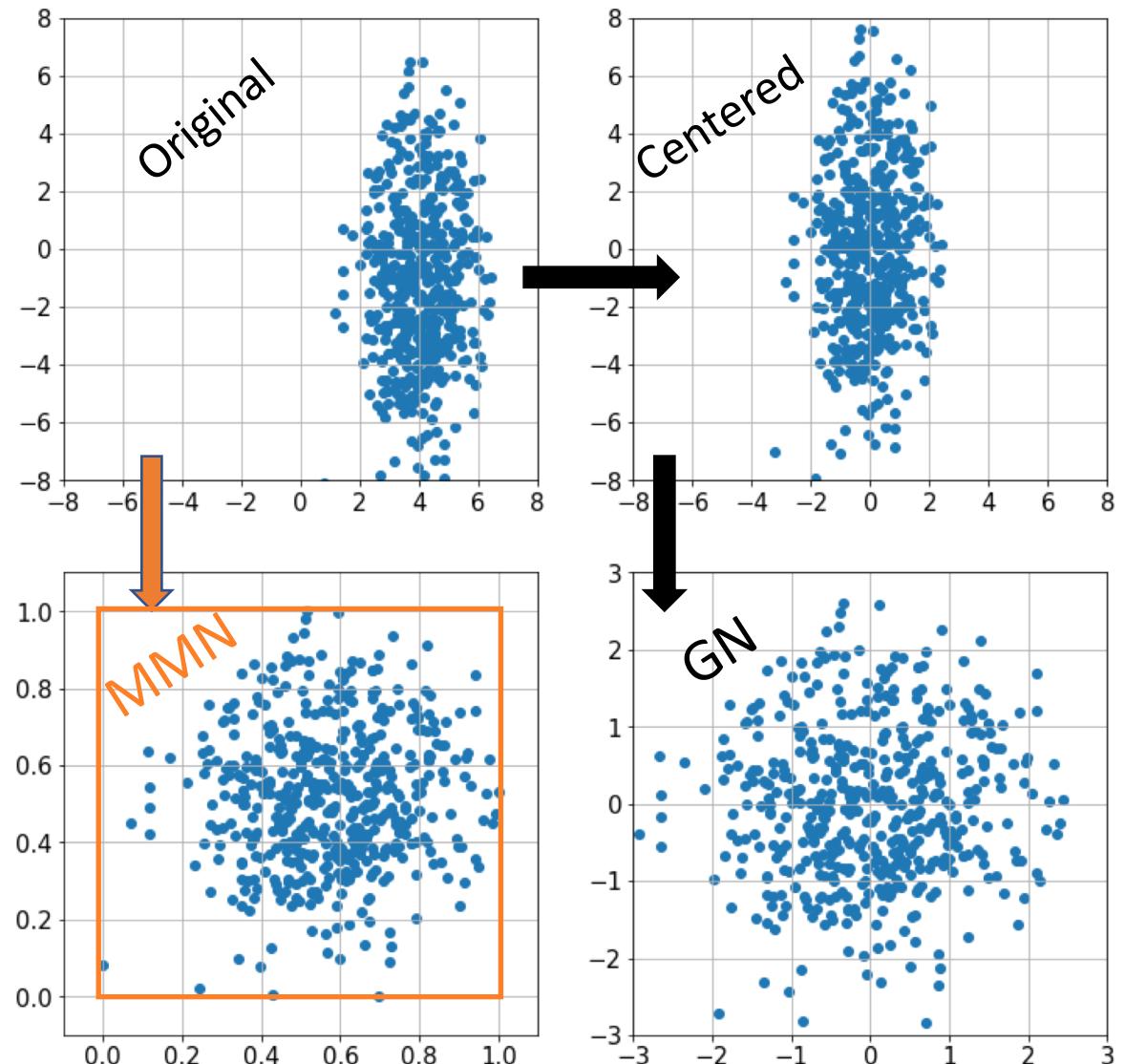
Essential preprocessing

$$X_{:,i} = \frac{X_{:,i} - \mu_i}{\sigma_i}$$

Gaussian normalization:
Each feature is a unit Gaussian

$$X_{:,i} = \frac{X_{:,i} - m_i}{M_i - m_i}$$

Minmax normalization:
Each feature is in [0,1]



Feature dependency and scaling

We want to make features independent

Why? Think dropout

$$\text{Cov}_{f \times f}(\mathbf{X}) = \frac{1}{n-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$$

$$\tilde{\mathbf{X}} = \mathbf{X} - \boldsymbol{\mu}$$

We want to make this diagonal

In fact, we want to make this Identity
to make features have same scale

--> *Whitening*

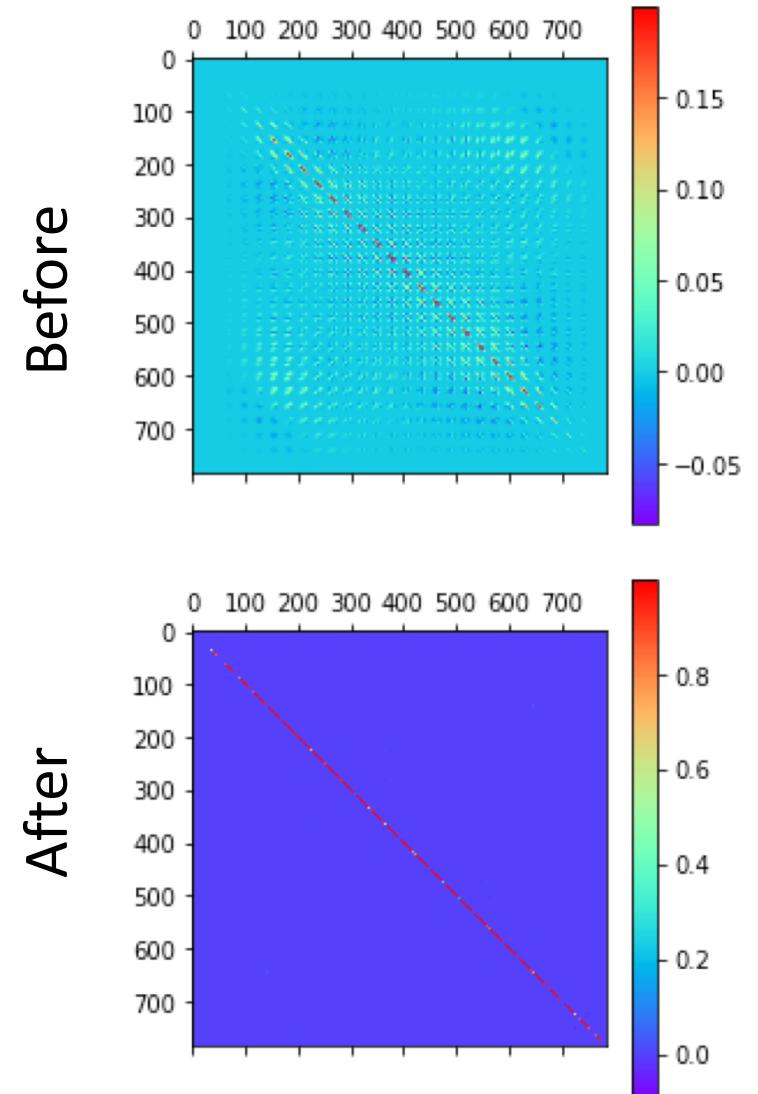
Zero Components Analysis (ZCA) Whitening

$$\begin{aligned} Z_{f \times f} &= \sqrt{n-1} (\tilde{X}^T \tilde{X})^{-\frac{1}{2}} \\ &= \sqrt{n-1} (Q D^{-\frac{1}{2}} Q^T) \end{aligned}$$

$$X_{\text{ZCA}} = \tilde{X} Z$$

Then...

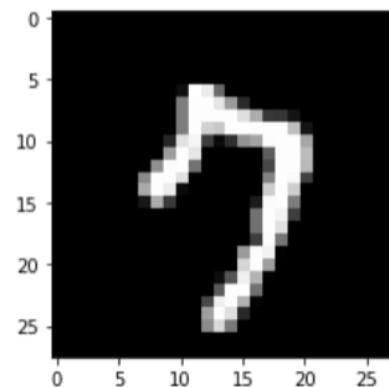
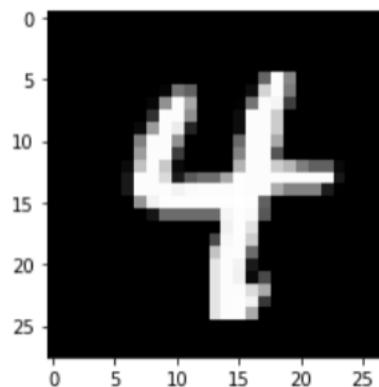
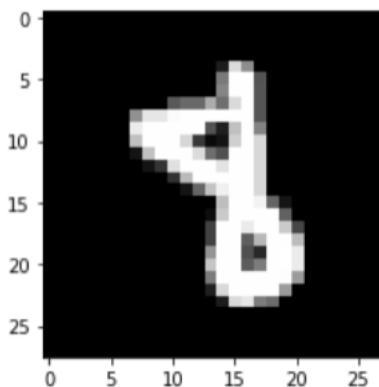
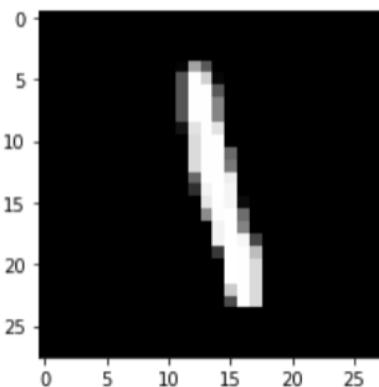
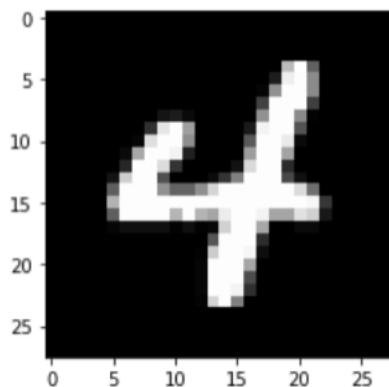
$$\frac{1}{n-1} X_{\text{ZCA}}^T X_{\text{ZCA}} = I$$



Dimensionality Reduction

Inputs can have a LOT of features

Eg: MNIST has 784 features, but we probably don't need the corner pixels



Simplify NN architecture by keeping only the essential input features

Principal Component Analysis (PCA)

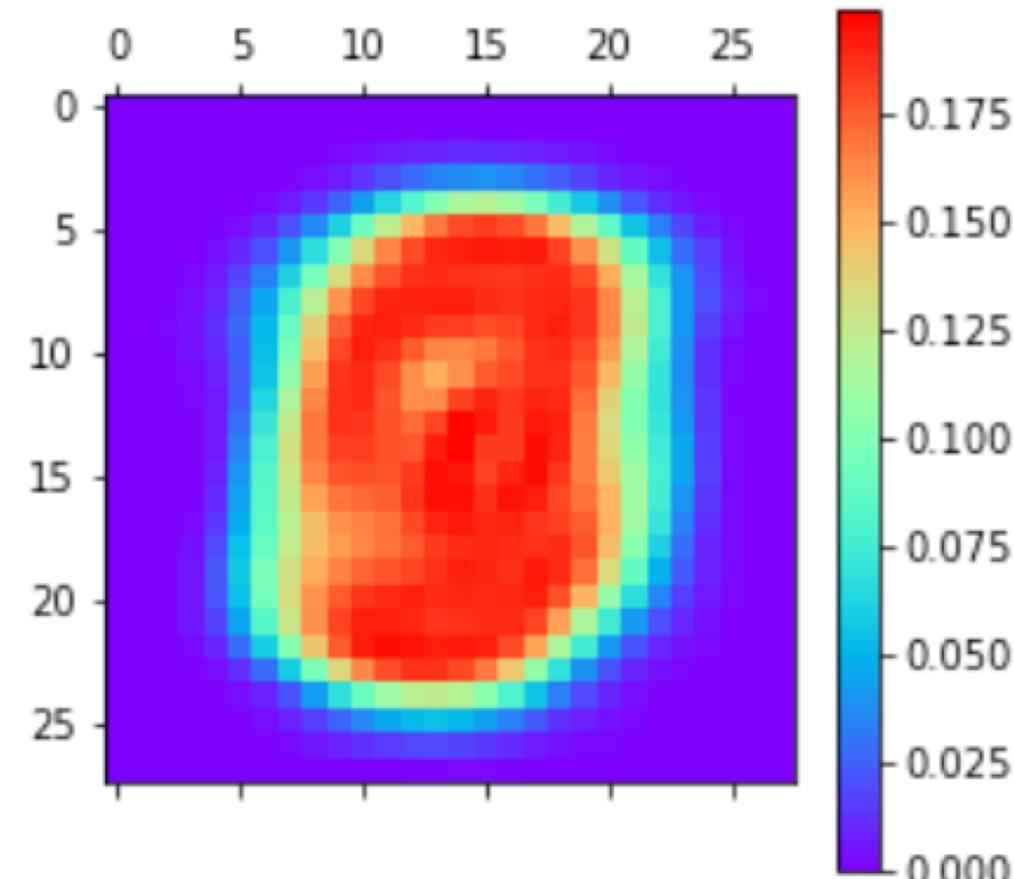
*Keep features with maximum variance
=> features which change the most*

Assume
sorted

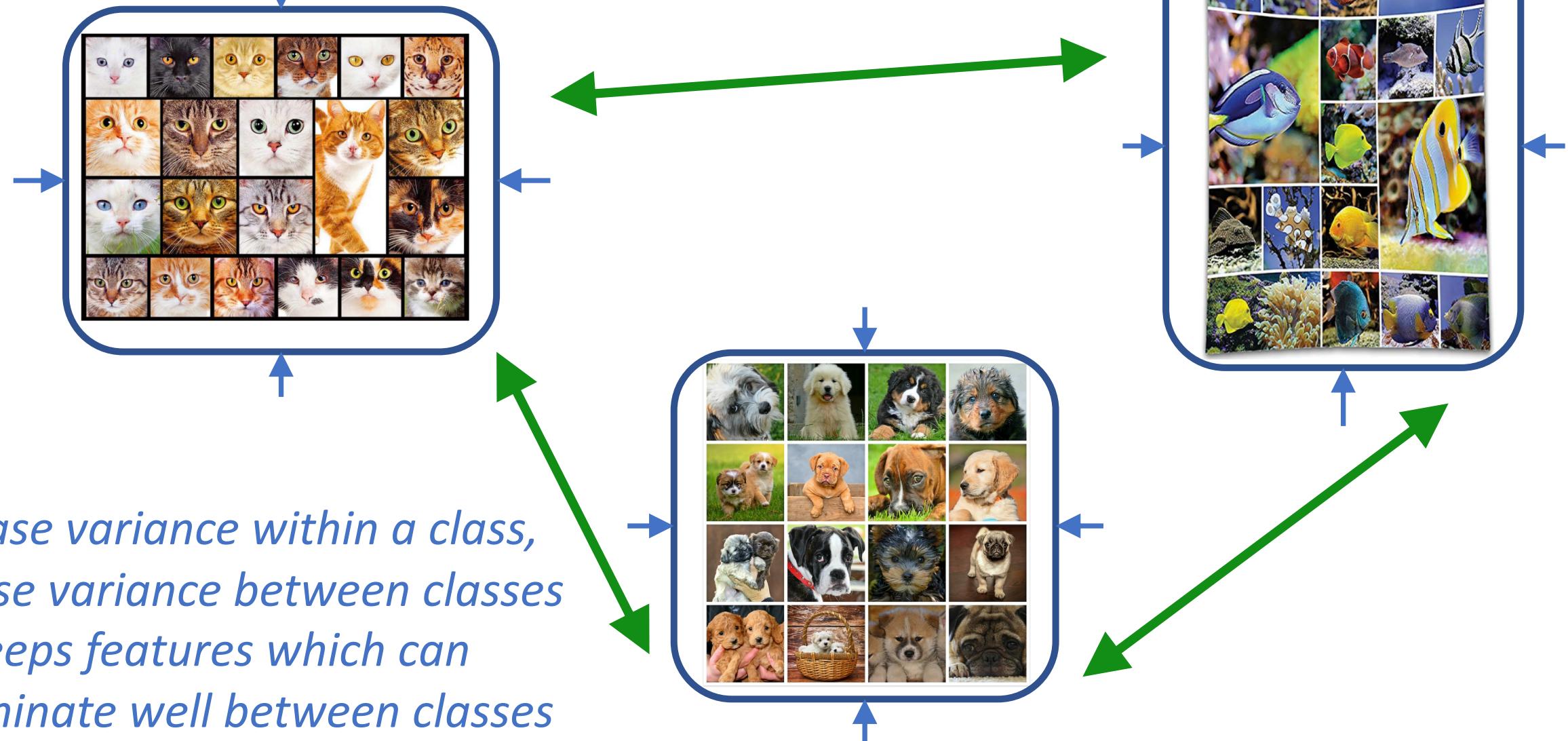
$$U_{f \times f} \Sigma_{f \times f} V_{f \times f} = \text{SVD} \{ \text{Cov}(X) \}$$

$$X_{\text{PCA}} = X_{n \times f'} U_{f \times f'}^{[:, :f']}$$

$f' < f$



Linear Discriminant Analysis (LDA)



Linear Discriminant Analysis (LDA)

Intra (within)
class scatter

Minimize

$$S_{\text{intra}} = \sum_{c=1}^C \left(\tilde{\mathbf{X}}_c - \boldsymbol{\mu}_c \right)^T \left(\tilde{\mathbf{X}}_c - \boldsymbol{\mu}_c \right)$$

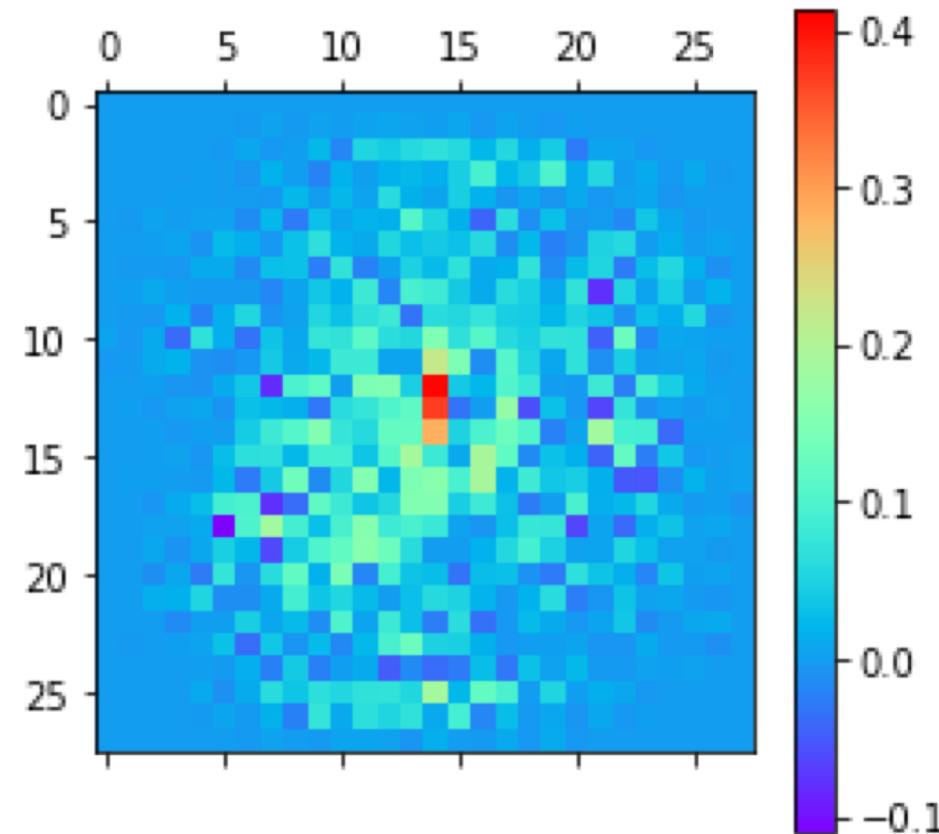
Inter (between)
class scatter

Maximize

$$S_{\text{inter}} = \sum_{c=1}^C N_c (\boldsymbol{\mu}_c - \boldsymbol{\mu}) (\boldsymbol{\mu}_c - \boldsymbol{\mu})^T$$

$$\mathbf{U}\Sigma\mathbf{V} = \text{SVD} \left\{ S_{\text{intra}}^{-1} S_{\text{inter}} \right\}$$

$$\mathbf{X}_{\text{LDA}} = \frac{\mathbf{X}}{n \times f} \quad \mathbf{U}[:, :f'] \quad f \times f'$$

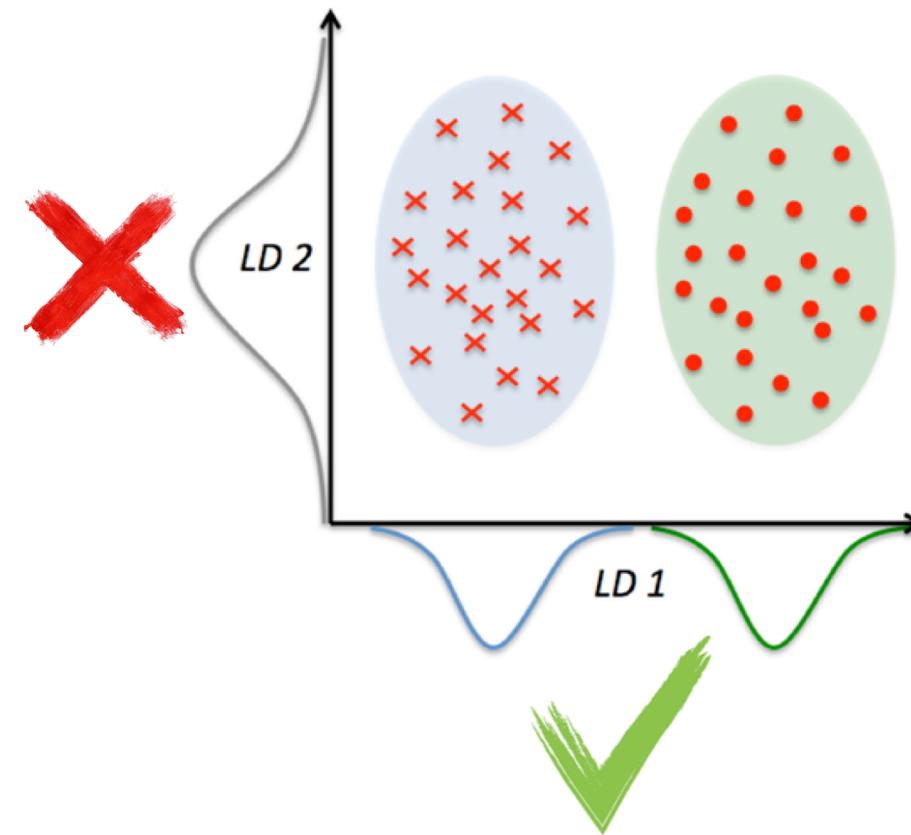
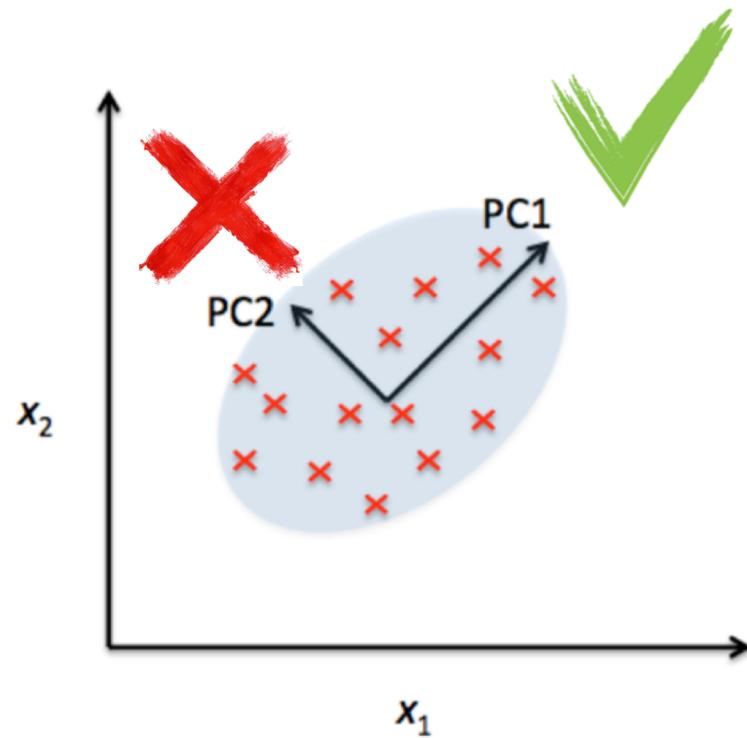


C: Number of classes

N_c : Number of samples in class c

$\boldsymbol{\mu}_c$: Feature mean of class c

PCA vs LDA



*PCA is unsupervised - just looks at data
LDA is supervised, looks at classes*

Remember...

*All statistics should be completed on training data only
Then apply them to validation and test data*

Example:

$$\mathbf{X}_{\text{train}} \rightarrow \mathbf{U}_{\text{train}}$$

$$\mathbf{U}_{\text{train}} \rightarrow \mathbf{X}_{\text{train}_{\text{PCA}}}$$

$$\mathbf{U}_{\text{train}} \rightarrow \mathbf{X}_{\text{val}_{\text{PCA}}}$$

$$\mathbf{U}_{\text{train}} \rightarrow \mathbf{X}_{\text{test}_{\text{PCA}}}$$

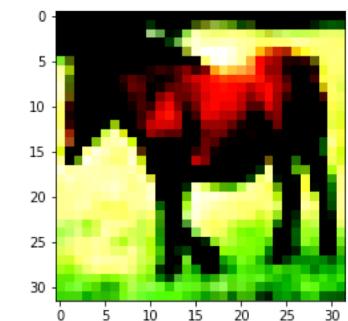
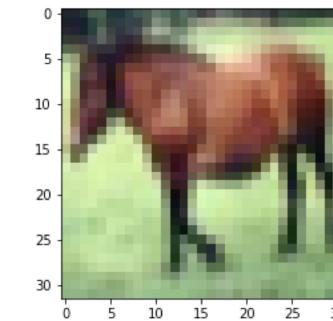
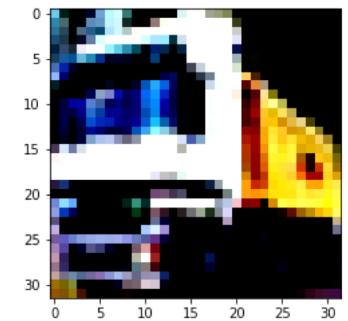
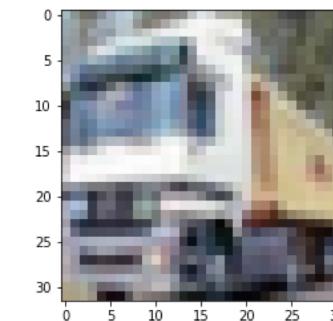
Global Contrast Normalization (GCN)

Increase contrast (standard deviation of pixels) of each image, one at a time

Different from other normalizations which operate on all training images together

\mathbf{x} : Single image

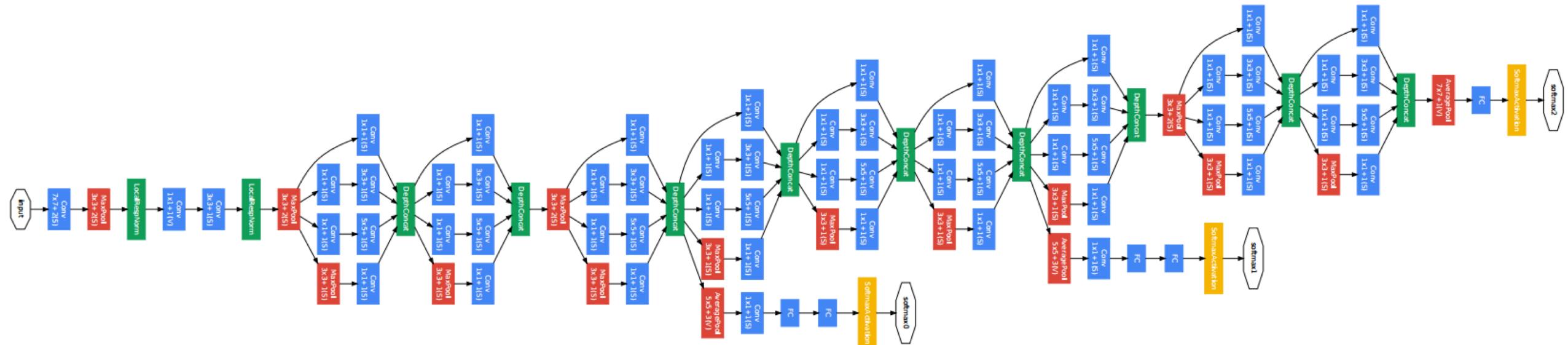
$$x_{\text{pixel}} = \frac{x_{\text{pixel}} - \mu(\text{all pixels in } \mathbf{x})}{\sigma(\text{all pixels in } \mathbf{x})}$$



Batch normalization

Input normalization scales the input features. But what about internal activations? Those are also input features to the layers after them

Internal activations get shifted to crazy values in deep networks



Batch normalization

$$\mu = \frac{1}{M} \sum_{i=1}^M x_i$$

Re-normalize internal values using a minibatch, then apply some trainable $(\mu, \sigma) = (\beta, \gamma)$ to them

x can be any value inside any layer of the network

$$\sigma^2 = \frac{1}{M} \sum_{i=1}^M (x_i - \mu)^2$$

$$\hat{x}_i = \frac{x_i - \mu}{\sigma + \epsilon}$$

Normalize

ϵ is just a small number to prevent division by 0. Can be machine epsilon

$$x_i \leftarrow \gamma \hat{x}_i + \beta$$

Scale the normalized values and substitute

Where to apply Batch Normalization?

Option 1: Before activation (at s)

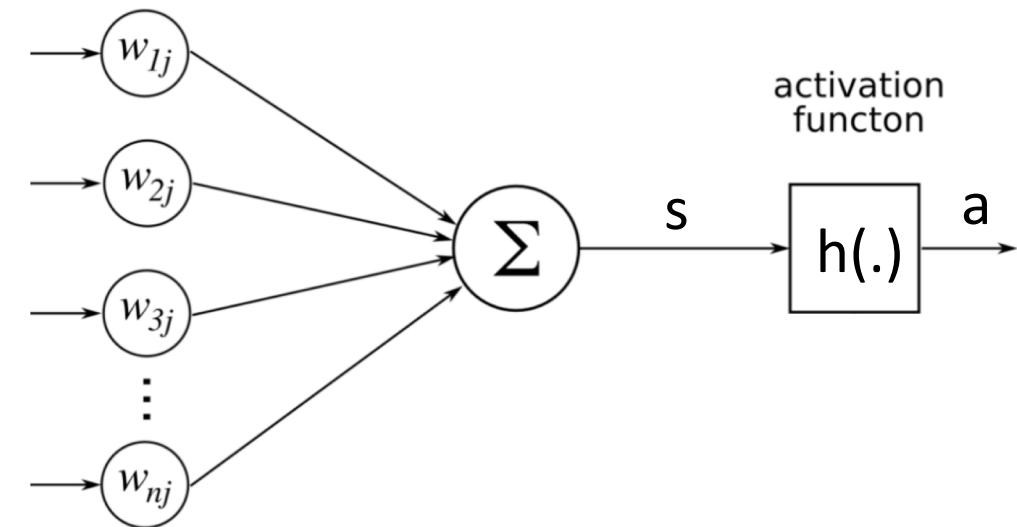
Option 2: After activation (at a)

Original paper recommends before activation.

This means scaling the input to ReLU, then ReLU can do its job by discarding anything ≤ 0 .

Scaling after ReLU doesn't make sense for β since it can be absorbed by the parameters of the **next** layer.

But... People still argue over this. I have done both and saw almost no difference.



$$\begin{aligned}s &= \mathbf{W}(\gamma \mathbf{a} + \beta) + \mathbf{b} \\&= (\mathbf{W}\gamma)\mathbf{a} + (\mathbf{W}\beta + \mathbf{b}) \\&= \tilde{\mathbf{W}}\mathbf{a} + \tilde{\mathbf{b}}\end{aligned}$$

Outline

- Activation functions
- Cost function
- Optimizers
- Regularization
- Parameter initialization
- Normalization
- **Data handling**
- Hyperparameter selection

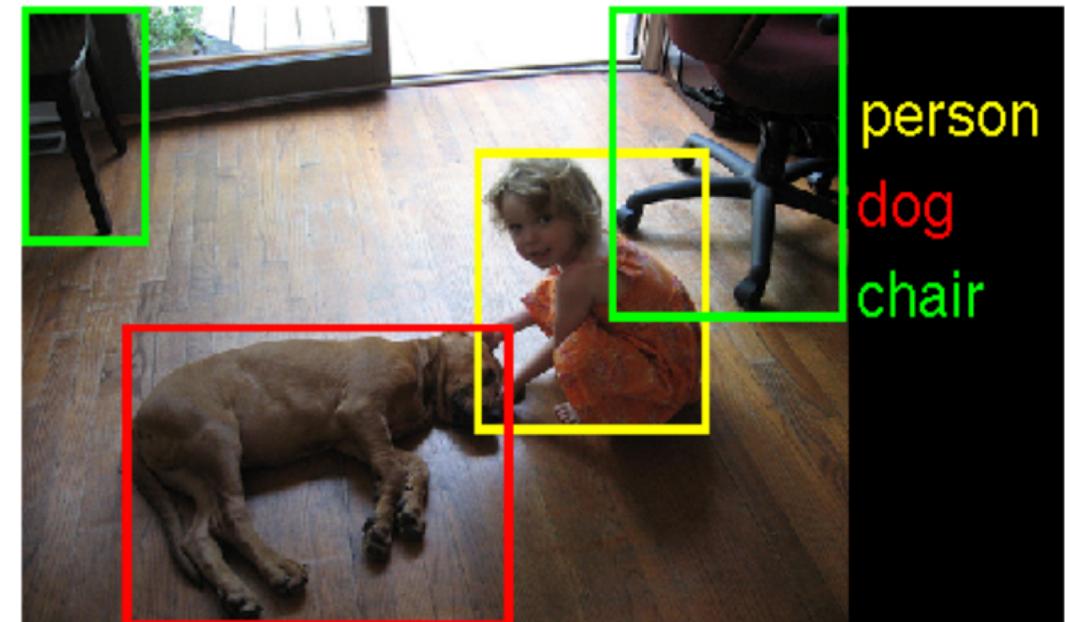
Data handling

- ✓ Collection
- ✓ Contamination
- ✓ Public datasets
- ✓ Synthetic data
- ✓ Augmentation



Collection and labeling

- Collected from real world sources - internet databases, surveys, paying people, public records, hospitals, etc...
- Labeled manually, usually crowd-sourced (Mechanical Turk)



Examples from Imagenet
<http://www.image-net.org>

Contamination



HW1: Human vs Computer, Binary
Did anyone cheat by using numpy? ;-)

Sometimes ground truth
labels can be unbelievable!

Ground truth	Predicted
6	5
5	8
9	6
6	1
1	4
4	9
9	5
1	1
9	4
5	3
6	0
4	9
3	5
0	6
9	4
5	3
6	1
4	9
3	5
5	3
6	1
2	7
1	7
7	1
6	1
3	5
8	3
7	1
0	7
2	7
4	9
9	5
6	6

Neural Networks are Lazy



Train



Test
Network output: Cat

For the NN, green background = cat

Missing Entries

Example: Adult Dataset

<https://archive.ics.uci.edu/ml/datasets/adult>

Features:

- Age
- Working class
- Education
- Marital status
- Occupation
- Race
- Sex
- Capital gain
- Hours per week
- Native country

Label: Income level

Missing at Random: Remove the sample

Missing not at Random: Removing sample may produce bias

Example: People with low education level don't reveal

Other ways:

Fill missing numerical values with mean or median. Eg: Age = 40

Fill missing categorical values with mode. Eg: Education = College

Small Dataset Size

Features (image
of breast cells):

- Radius
- Perimeter
- Smoothness
- Compactness
- Symmetry

Labels:

- Malignant
- Benign

Example: Wisconsin Breast Cancer Dataset

[https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+\(original\)](https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(original))

[https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic))

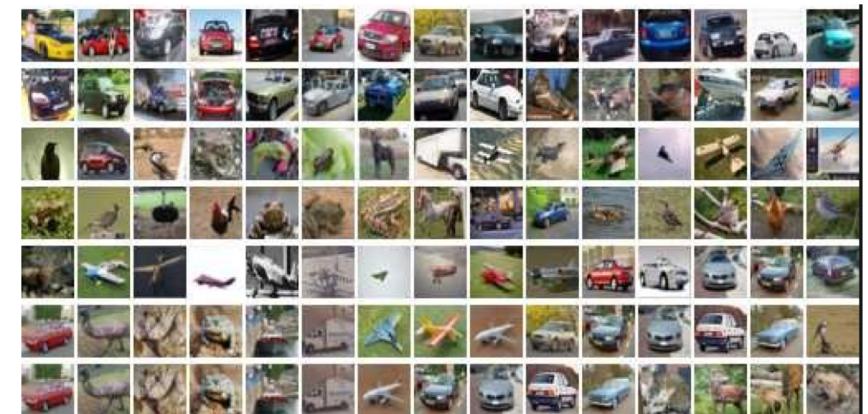
Original dataset: 699 samples with missing values

Modified dataset: 569 samples

*Such small datasets generally not suited
for neural networks due to overfitting*

Commonly used Image Datasets

- MNIST (MLP / CNN):
 - 28x28 images, 10 classes
 - Initial benchmark
 - If you're writing papers, don't just do MNIST!
 - SOTA testacc: >99%
 - **Harder Variation: Fashion MNIST**
- CIFAR-10, -100 (CNN):
 - 32x32x3 images (RGB), 10 or 100 classes
 - Widely used benchmark
 - SOTA testacc: ~97%, ~84%



Commonly used Image Datasets

- Imagenet (CNN):
 - Overall database has >14M images
 - ILSVRC challenge has >1M images, 224x224x3, 1000 classes
 - Widely used benchmark, don't attempt without enough computing power!
 - Alexnet top-1, top-5 testacc: 62.5%, 83% (not SOTA, but very famous!)
 - **Simpler variation: Tiny Imagenet**
- Microsoft Coco (CNN / Encoder-Decoder):
 - 330k images, 80 categories
 - Segmentation, object detection
- Street View House Numbers (CNN / MLP):
 - Think MNIST, but more real-world, harder, bigger
 - SOTA testacc > 98%



Other Commonly used Datasets

- TIMIT (Preprocessing + MLP + Decoder):
 - Automatic Speech recognition (Guest lecture coming up on 2/27)
 - Speaker and phoneme classification
- Reuters (MLP):
 - Classify articles into news categories
 - Multiple / tree structured labels
- IMDB Reviews:
 - Natural language processing
- YouTube-8M:
 - Annotate videos

Source for Datasets

- Kaggle: <https://www.kaggle.com/datasets>
- UCI ML repository: <https://archive.ics.uci.edu/ml/datasets.html>
- Google: <https://toolbox.google.com/datasetsearch>
- Amazon Web Services: <https://registry.opendata.aws/>

Synthetic Datasets

*Data that is artificially manufactured (using algorithms),
rather than generated by real-world events*

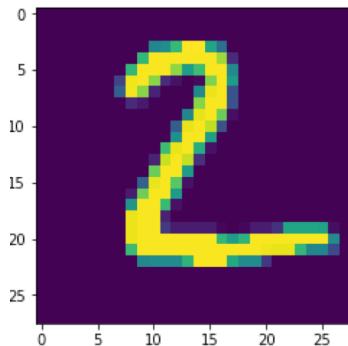
- Tunable algorithms to mimic real-world as required
- Cheap to produce large amounts of data
- We created data on Morse code symbols with tunable classification difficulty:
<https://github.com/usc-hal/morse-dataset>
- Can be used to augment real-world data

Data Augmentation

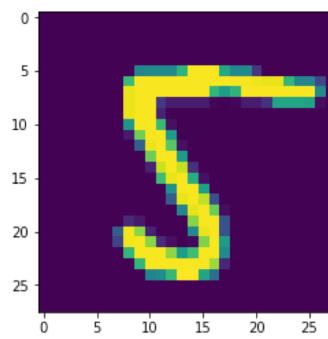
*More training results in overfitting
The solution: Create more data!*

Eg: More MNIST images can simulate more different ways of writing digits
If the network sees more, it can learn more

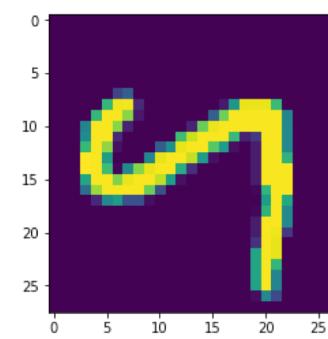
Data Augmentation Examples for Images



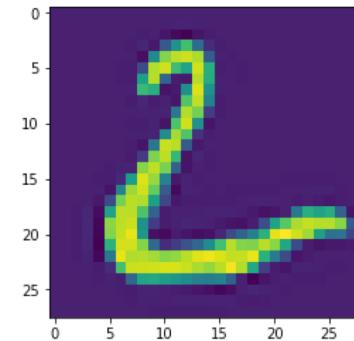
Original



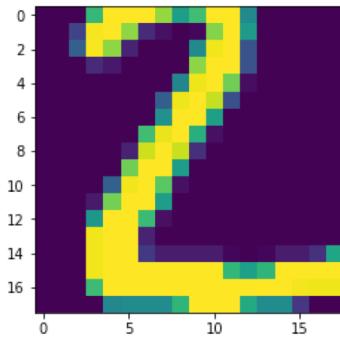
Flip Top-Bottom
Do NOT do for digits



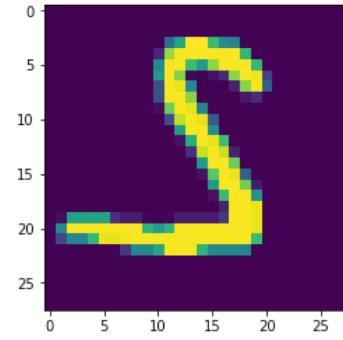
Transpose



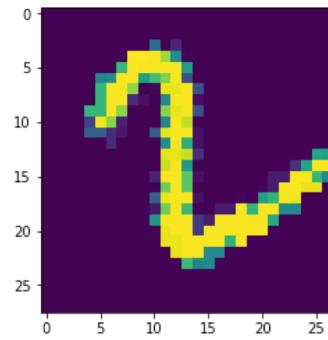
Elastic
Transform



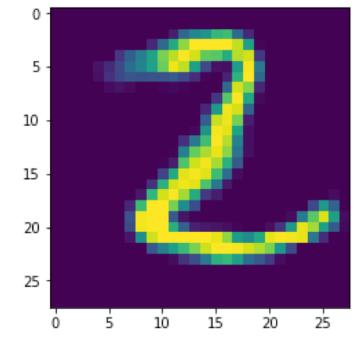
Cropped



Flip Left-Right



Rotate 30°



Simard, Steinkraus and Platt, "Best Practices for Convolutional Neural Networks applied to Visual Document Analysis", in Proc. Int. Conf. on Document Analysis and Recognition, 2003.

Data Augmentation in Python

I use the Pillow Imaging library

```
>> pip install Pillow

from PIL import Image
im = Image.fromarray((x.reshape(28,28)*255).astype('uint8'), mode='L')
im.show()
im.transpose(Image.FLIP_LEFT_RIGHT)
im.transpose(Image.FLIP_TOP_BOTTOM)
im.transpose(Image.TRANSPOSE)
im.rotate(30)
im.crop(box=(5,5,23,23)) #(left,top,right,bottom)
```

Outline

- Activation functions
- Cost function
- Optimizers
- Regularization
- Parameter initialization
- Normalization
- Data handling
- Hyperparameter selection

What is a Hyperparameter?

*Anything which the network does NOT learn,
i.e. its value is adjusted by the user*

Continuous Examples:

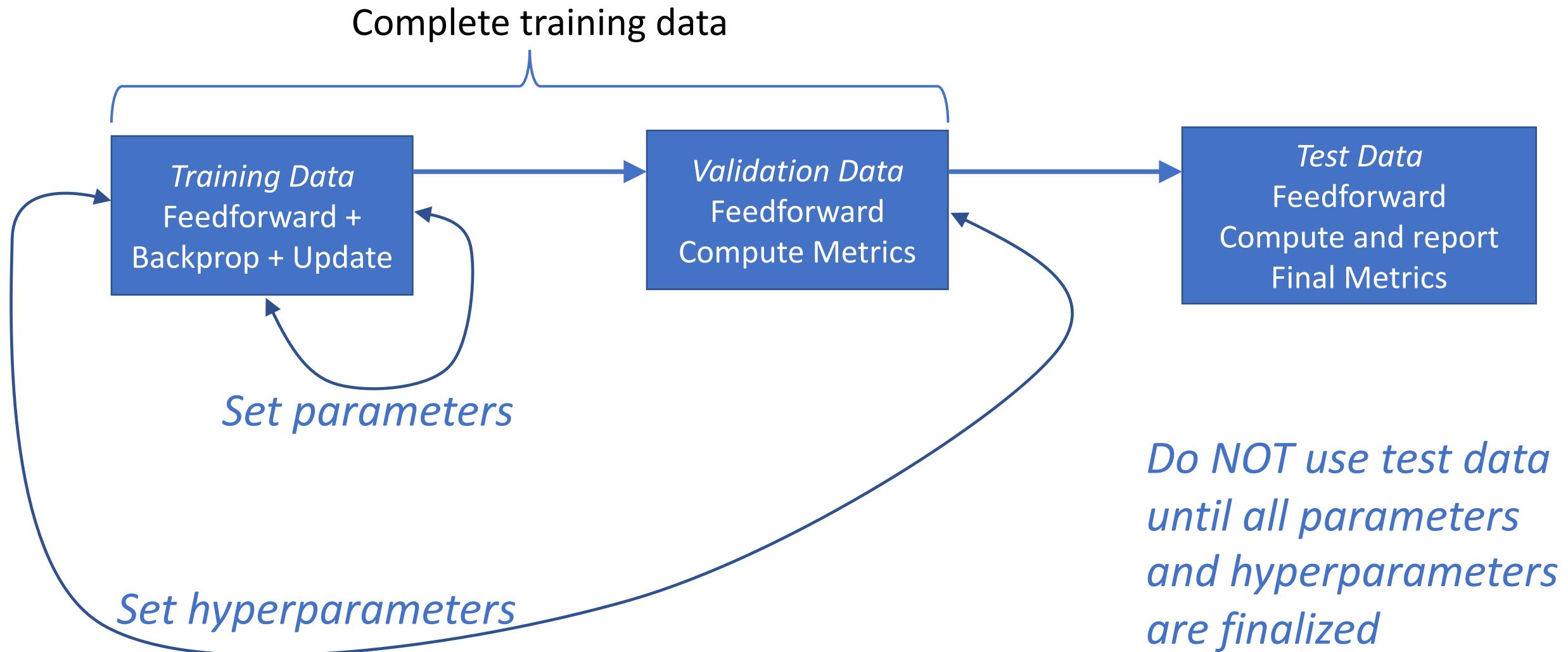
- Learning rate η
- Momentum α
- Regularization λ



Discrete Examples:

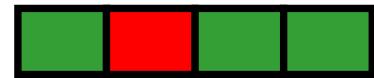
- What optimizer: *SGD, Adam, RMSprop, ...?*
- What regularization: *L2, L1, both, ...?*
- What initialization: *Glorot, He, normal, uniform, ...?*
- Batch size M: *32, 64, 128, ...?*
- Number of MLP hidden layers: *1, 2, 3, ...?*
- How many neurons in each layer?
- What kinds of data augmentation?

Using data properly



Cross-Validation

- Divide the dataset into k parts.
- Use all parts except i for training, then test on part i .
- Repeat for all i from 1 to k .
- Average all k test metrics to get final test metric.



*Useful when dataset is small
Generally not needed for NN problems*

Hyperparameter Search Strategies

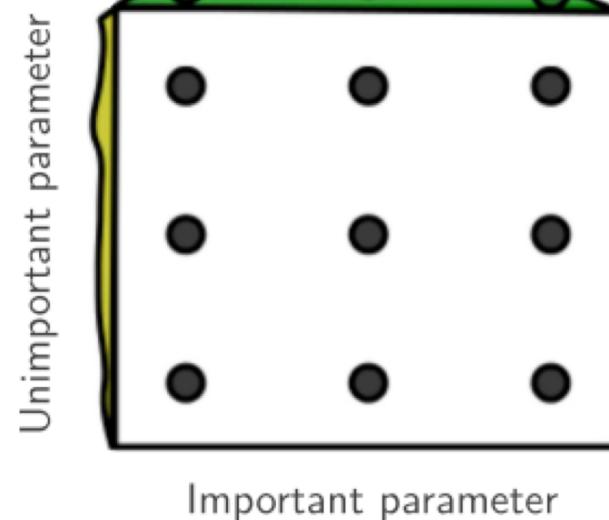
Grid search for (η, α) :

$(0.1, 0.5), (0.1, 0.9), (0.1, 0.99),$
 $(0.01, 0.5), (0.01, 0.9), (0.01, 0.99),$
 $(0.001, 0.5), (0.001, 0.9), (0.001, 0.99)$

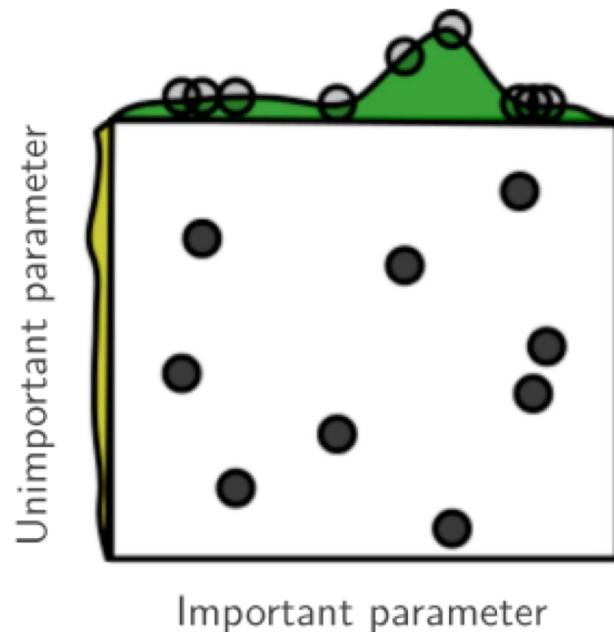
Random search for (η, α) :

$(0.07, 0.68), (0.002, 0.94), (0.008, 0.99),$
 $(0.08, 0.88), (0.005, 0.62), (0.09, 0.75),$
 $(0.14, 0.81), (0.006, 0.76), (0.01, 0.8)$

Grid Layout



Random Layout



*Random search samples more values
of the important hyperparameter*

Pic courtesy: J. Bergstra, Y. Bengio, "Random Search for Hyperparameter Optimization," Journal of Machine Learning Research, vol. 13, pp. 281–305, 2012

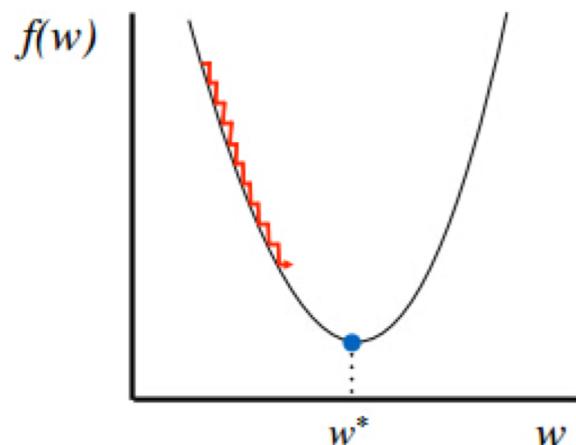
Hyperparameter Optimization Algorithm : Tree-Structured Parzen Estimator (TPE)

Potential project topic

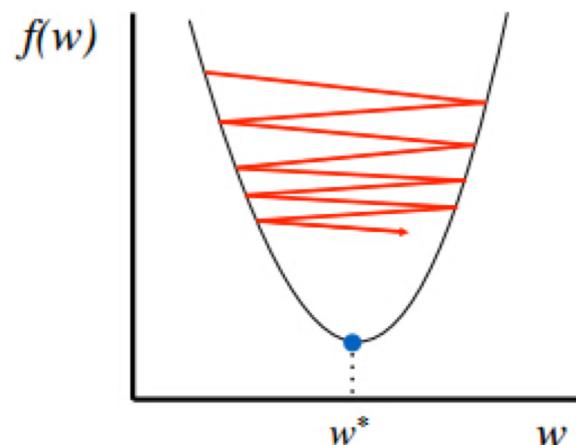
1. Decide initial probability distribution for each hyperparameter
→ Example: $\eta \sim \text{logUniform}(10^{-4}, 1)$; $\lambda \sim \text{logUniform}(10^{-6}, 10^{-3})$
2. Set some performance threshold
→ Example: MNIST validation accuracy after 10 epochs = 96%
3. Sample multiple *hyp_sets* from the distributions
4. Train, and compute validation metrics for each *hyp_set*
5. Update distribution and threshold:
 - If performance of a *hyp_set* is worse than threshold: Discard
 - If performance of a *hyp_set* is better than threshold: Change distribution accordingly
6. Repeat 3-5

Learning rate schedules

When using simple SGD, start with a big learning rate, then decrease it for smoother convergence



Too small: converge very slowly



Too big: overshoot and even diverge

Learning rate schedules

$$\eta_{t+1} = \frac{\eta_0}{1 + kt}$$

Initial η →
Epochs →

Fractional decay

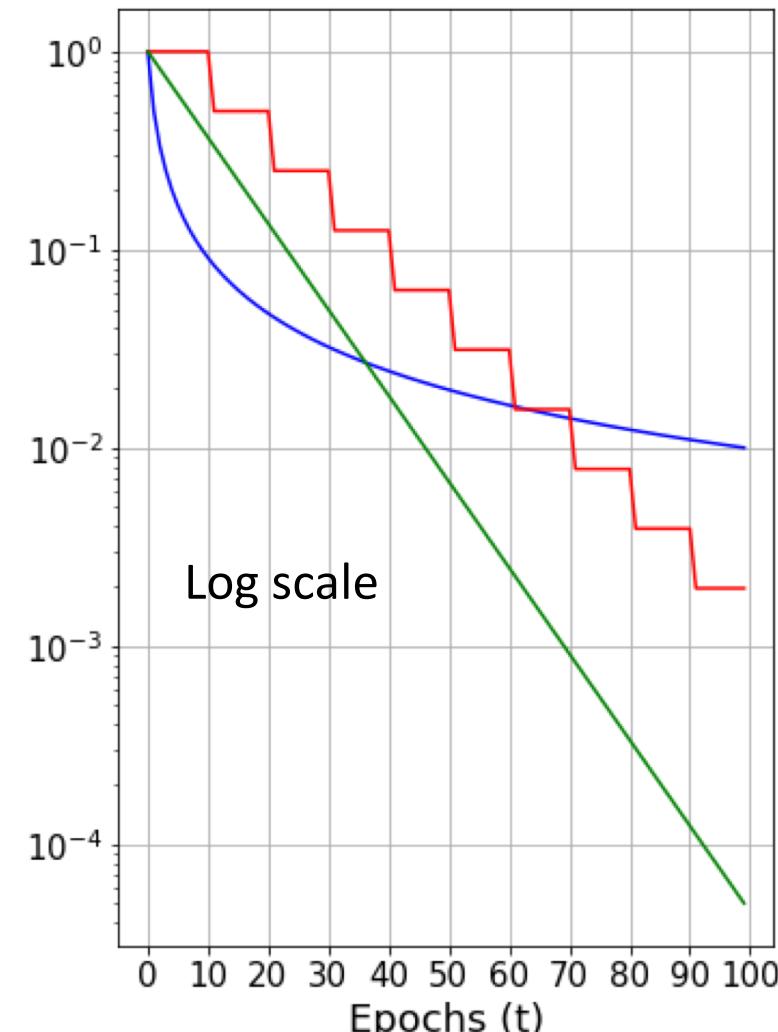
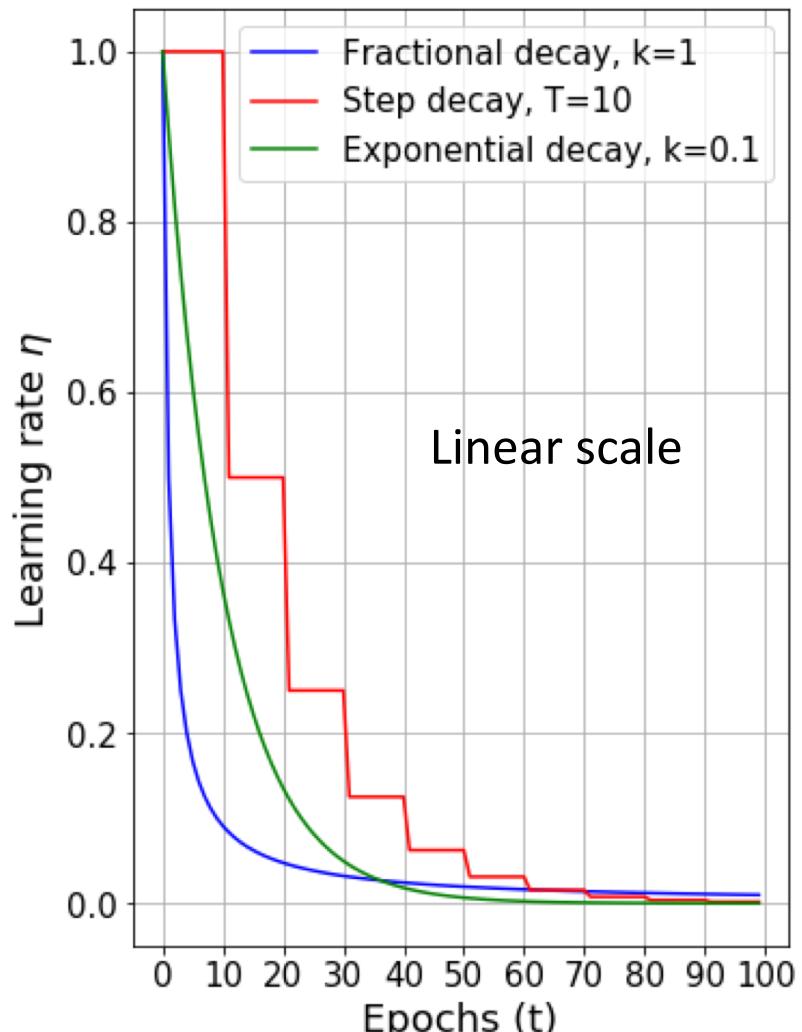
Used in Keras

$$\eta_{T+1} = \eta_0 \left(\frac{1}{2}\right)^{\lfloor t/T \rfloor}$$

Step decay

$$\eta_{t+1} = \eta_0 e^{(-kt)}$$

Exponential decay



Other Learning rate strategies

Potential project topic

Warmup: Increase η for a few epochs at the beginning. Works well for large batch sizes.

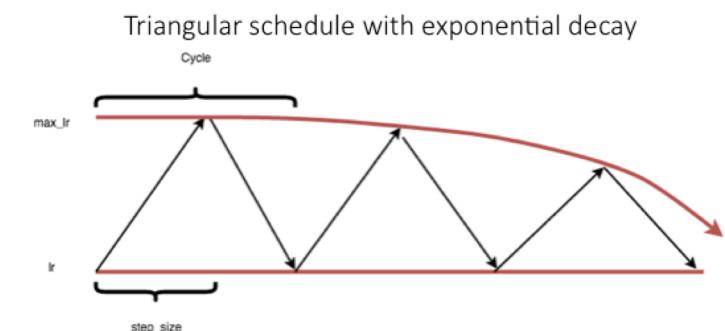
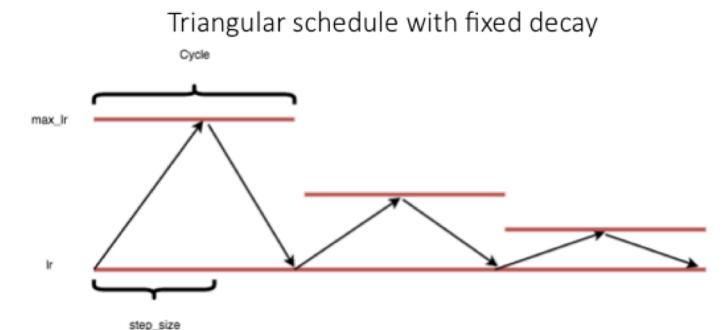
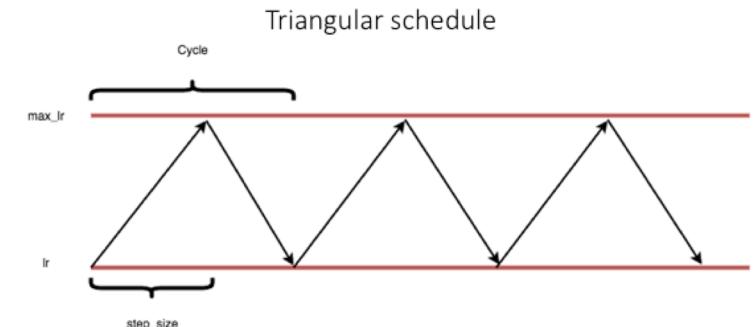
Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, Kaiming He. "Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour". arXiv:1706.02677

Cosine Scheduling

I. Loshchilov, F. Hutter, "SGDR: Stochastic Gradient Descent with Warm Restarts", Proc. ICLR 2017.

Triangular Scheduling

L. N. Smith, "Cyclical Learning Rates for Training Neural Networks", arXiv:1506.01186
Pic courtesy <https://www.jeremyjordan.me/nn-learning-rate/>



That's all
Thanks!

