

Have to combine disparate information sources.

Alpha Go Zero: 100-0 vs Alpha Go. Learned entirely from self play
 • Always has an opponent at just the right level

Algorithms matter way more than data/compute... principled approach

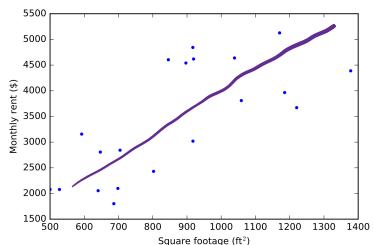
- Label smoothing - Don't treat data from Humans as 100% accurate

History of NN's

- Rosenblatt : Perceptron, lines > 0 , else doesn't fire
 - Is a linear classifier
 - Multi-layer perceptrons achieved non-linear behaviors
 - Needed back prop to train
- Fully connected , convolutional Neural Network
 - Why revolution post 2010? Datasets and compute
 - ImageNet
 - 14 million Images
 - Bounding Boxes - localization
 - Alex Net - 1st to get 15% error on ImageNet, kicked off CNN revolution

Basics of ML : LinearReg

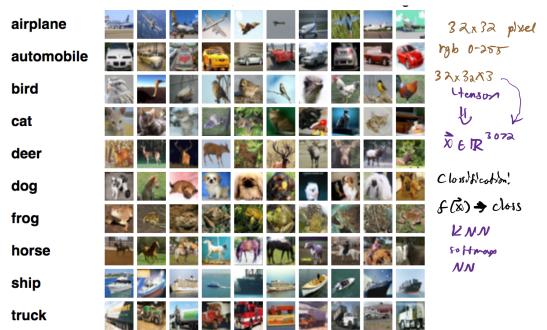
- Focus on supervised classification problems
- CIFAR-10 dataset used a lot
- Linear reg example



input: sq ft
output: rent

Model? reasonably linear

Read Ch. 5 - 5.5 2,3



• Min J(a,b)

$$y = a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n$$

parameters: choose a, b for best fit

- Could have done a polynomial: $y = b + a_{100}x^{100} + a_{99}x^{99} + \dots + a_1x$

- Highly non-linear and overfit
- Degree of Freedom ~ # of parameters
- Choose using cross validation, test set
- NN is a Universal Function Approx: Can emulate any function that exists
- Why sq' error?: Derivative calc is easier

Lec 2: Linear Reg

1/5

Linear Reg Example

$$\text{Model: } y = ax + b = \theta^T x \quad \theta = \begin{bmatrix} a \\ b \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ x \end{bmatrix}$$

$$\text{Cost: } \mathcal{L}(\theta) = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2 \quad \hat{y} = \theta^T x$$

↑ scaling mag

Solve by setting derivatives / gradient $\frac{d\mathcal{L}}{d\theta} = 0$

Multivariate Derivatives / Gradients:

$$x \in \mathbb{R}^n, y \in \mathbb{R}$$

1) $\nabla_x y = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}$ - each number tells you Δx_i gives you Δy
 • number in linear reg. function is general
 $(\nabla_x y)_2$
 $\Delta y \text{ due to } x_2 \text{ is } \left\{ \frac{\partial y}{\partial x_2} \right\} \Delta x_2$

2) $f(x) = x^T A x, \nabla_x f(x)? \quad x \in \mathbb{R}^n \quad A \in \mathbb{R}^{n \times n}$

$$A = \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{bmatrix} \quad f(x) = \sum_i \sum_j a_{ij} x_i x_j$$

$$\sum_{i,j=1} \frac{\partial a_{ij} x_i^2}{\partial x_i} = 2 a_{ii} x_i, \sum_{i,j \neq i} \sum_{j=2} a_{ij} x_i x_j = \sum_{j=2} a_{jj} x_j \Rightarrow \frac{\partial f(x)}{\partial x_i} = 2 a_{ii} x_i + \sum_{j=2} a_{ij} x_j + \sum_{i=2} a_{ii} x_i$$

$$\sum_{j=1} a_{ij} x_j + \sum_{i=1} a_{ii} x_i \quad (Ax)_i \quad (A^T x)$$

$$\nabla_x f(x) = (A x + A^T x) = (A + A^T) x$$

$$A = A^T \quad \text{if } A \text{ is symmetric} = 2 A x$$

Denominator layout

$\mathbb{D} \xrightarrow{\text{denom} \times \text{num}}$

$\text{dim} = \text{denom}, \text{num}^T$

Numerator layout

$\mathbb{D} \xrightarrow{\text{num} \rightarrow \text{denom}}$

$\text{dim} = \text{num}, \text{denom}^T$

Least Squ. Solution:

$$\theta = (X^T X)^{-1} X^T y$$

$$\begin{aligned}\frac{\partial L}{\partial \theta} &= \frac{1}{2} [D - 2X^T y + (X^T X + X^T X) \theta] \\ &= -X^T y + X^T X \theta = 0 \\ \theta &= (X^T X)^{-1} X^T y\end{aligned}$$

For higher order:

$$x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \theta = \begin{bmatrix} b \\ a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix}$$

x - features vector, y is values

More data - more complex models

Lec 3

k-fold / MLE

110

θ - parameters

Hyper - choices prior to fitting

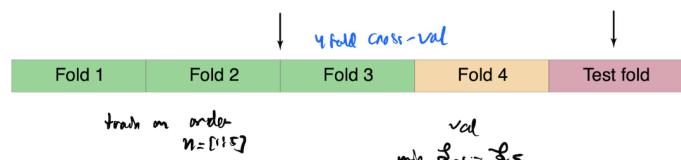
The standard for training, evaluating, and choosing models is to use different datasets for each step.

get θ

- Training data is data that is used to learn the parameters of your model.
- Validation data is data that is used to optimize the hyperparameters of your model. This avoids the potential of overfitting to nuances in the testing dataset. ~~fit~~ ~~hyperparams~~
- Testing data is data that is used to score your model. ~~- unseen until testing~~

k-fold cross validation

k-folds in train/val not test



Max Likelihood Est (MLE)

$$\mathcal{L} = \prod_{i=1}^N p(x_i, y_i | \theta)$$

$$\log \mathcal{L} = \sum \log p(x_i, y_i | \theta) = \sum \log [p(y_i) p(x_i | y_i)]$$

$$= N \log \underbrace{\frac{1}{\sqrt{2\pi}}}_{\text{equimprobable}} + \sum \log N(\mu_{y(i)}, \Sigma_{y(i)})$$

$$\begin{aligned}\log p(x^{(i)} | y^{(i)} = j) &= \log \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right) \\ &= -\frac{n}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \\ &= k_2 - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\end{aligned}$$

Hence, we arrive at:

$$\log \mathcal{L} = k + \sum_{i=1}^N -\frac{1}{2} \log |\Sigma_{y(i)}| - \frac{1}{2} (x^{(i)} - \mu_{y(i)})^T \Sigma_{y(i)}^{-1} (x^{(i)} - \mu_{y(i)})$$

• Plurality of k nearest (Euclidean Dist)

1. Chose metric: Euclidean $\|x^i - x^j\|$
2. # neighbors k
3. $d(x^{\text{new}}, x^i) \forall i$
4. $\{c_1, \dots, c_k\}$ smallest distances, x^{new} is most common class

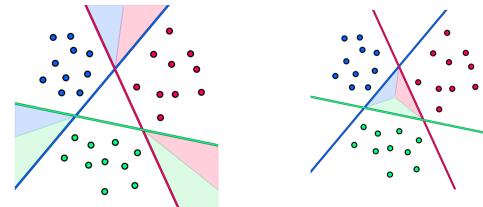
Pros: simple Cons: slow class, need every point in memory

Distances are far less intuitive in higher dim space

What's linear class doing?

Further distance away on positive side is better!

$$\max w_i^T x \quad \text{if } w_i = \frac{w}{\|w\|_2}$$



https://jermwatt.github.io/machine_learning_refined/notes/7_Linear_multiclass_classification/7_2_OvA.html

Soft max

$$a_i = w_i^T x \quad \frac{e^{a_i(x)}}{\sum_{j=1}^c e^{a_j(x)}}$$

$$\Rightarrow \Pr(y^i = i | x^i, \theta) = \text{softmax}_i(x^i)$$

Prob that x^i is in class j

Lec 5 - log softmax, Gradient Descent

Soft max Classification

$$\arg \max_{\theta} \prod_{i=1}^m p(x^i | \theta) p(y^i | x^i, \theta) = \arg \max_{\theta} \prod_{i=1}^m p(y^i | x^i, \theta)$$

equiprobable priors

Soft max log likelihood:

$$\arg \min_{\theta} \sum_{i=1}^m \left(\log \sum_{j=1}^c e^{a_j(x)} - a_{y(i)}(x^{(i)}) \right)$$

$$\arg \max_{\theta} \sum_{i=1}^m \log p(y^i | x^i, \theta) = \arg \max_{\theta} \sum_{i=1}^m \log \text{softmax}_{y(i)}(x^i)$$

$$= \arg \max_{\theta} \sum_{i=1}^m \log \left[\frac{e^{a_{y(i)}(x^i)}}{\sum_{j=1}^c e^{a_j(x^i)}} \right]$$

$$= \arg \max_{\theta} \sum_{i=1}^m \left[a_{y(i)}(x^i) - \log \sum_{j=1}^c e^{a_j(x^i)} \right] \frac{1}{m}$$

$$= \arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m \left[\log \sum_{j=1}^c e^{a_j(x^i)} - a_{y(i)}(x^i) \right]$$

make a loss

Gradient Descent:

• 1st order approx $\nabla_x f \cdot \Delta x \approx f(x + \Delta x)$
 linearization

• Step size:



Remember, step size is using linear approximation

- Batch vs Mini (Stochastic) -
 - Batch is all at once (all in memory)
 - Mini batch GD - Use a small subset 1-10k items
 - Stochastic - one example and step

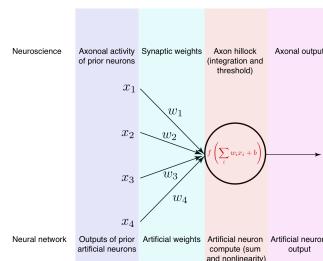
Lec 6 : NN's : AE, Activation etc

1124

Neural Networks:

• Bio neuron

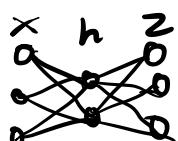
- axon, axon hillock etc:



of layers: hidden + output \rightarrow n hidden + 1

If $d(m(h)) \ll d(m(x))$: Low rank representation of data

Auto encoder: lower dim represent of x



$$\mathcal{L} = \|z - x\|^2$$

PCA is linear, auto encoder doesn't have to be

Activation Functions (Pros + cons):

Sigmoid:

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

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

Pros:

- Around $x = 0$, the unit behaves linearly.
- It is differentiable everywhere.

Cons:

- At extremes, the unit saturates and thus has zero gradient. This results in no learning with gradient descent.
- The sigmoid unit is not zero-centered; rather its outputs are centered at 0.5.

Consider $f(w) = \sigma(w^T x + b)$. Defining $z = w^T x + b$, the derivative with respect to w , the parameters, is:

$$\frac{\partial f(w)}{\partial w} = \sigma(z)(1 - \sigma(z))x_{h_1} = \frac{25}{22} \frac{26}{22} \frac{26}{22} \text{ for } \sigma(z)$$

If $x \geq 0$ (e.g., if the input units all had a sigmoidal output), then the gradient has all positive entries. Let's say we had some gradient, $\frac{\partial \mathcal{L}}{\partial w}$, which can be positive or negative. Then $\frac{\partial \mathcal{L}}{\partial w}$ will have all positive or negative entries.

This can result in zig-zagging during gradient descent.



Hyperbolic Tangent:

$$\tanh(x) = 2\sigma(x) - 1$$

$$\tanh'(x) = 1 - \tanh^2(x)$$

Pros:

- Around $x = 0$, the unit behaves linearly.
- It is differentiable everywhere.
- It is zero-centered.

Cons:

- Like the sigmoid unit, when a unit saturates, i.e., its values grow larger or smaller, the unit saturates and no additional learning occurs.

ReLU

$$\text{ReLU}(x) = \max(0, x)$$

$$\text{ReLU}'(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$

Pros:

- In practice, learning with the ReLU unit converges faster than sigmoid and tanh.
- When a unit is active, it behaves as a linear unit.
- The derivative at all points, except $x = 0$, is 0 or 1. When $x > 0$, the gradients are large, and not scaled by second order effects.
- There is no saturation if $x > 0$.

Cons:

- $\text{ReLU}(x)$, like sigmoid, is not zero-centered.
- $\text{ReLU}(x)$ is not differentiable at $x = 0$. However, in practice, this is not a large issue. A heuristic when evaluating $\frac{d\text{ReLU}(x)}{dx} \Big|_{x=0}$ is to return the left derivative (0) or the right derivative (1); this is reasonable given digital computation is subject to numerical error.
- Learning does not happen for examples that have zero activation. This can be fixed by e.g., using a leaky ReLU or maxout unit.

ReLU Variants:

Softplus: $\text{f}_n(x) = \log(1 + \exp(x))$

Leaky ReLU: $\text{max}(\alpha x, x)$

Exp ReLU: $\max(\alpha(\exp(x)-1), x)$

Output activation / Loss

$\sigma(z)$

$$MSE = \frac{1}{n} \sum_i [y_i - \sigma(z_i)]^2$$

$$\text{Cross Entropy (CE)} = - \sum_i y_i \log \sigma(z_i) + (1-y_i) \log(1-\sigma(z_i))$$

$$\frac{\partial \text{MSE}}{\partial z} = \sqrt{\text{learning steps}}$$

$$\frac{\partial \text{CE}}{\partial z} = \sqrt{\text{learning steps near output layer}}$$

Lec 7: Back Prop, Gates, NN Backprop Layer

1126

Backprop is Chain Rule, local gradient with upstream derivatives

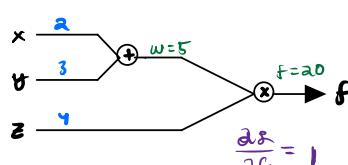
If you know local gradients, you can backprop anything

$$f(x,y,z) = (x+y)z \quad u = x+y \quad \lambda = 1-f(u)$$

Due to denominator - we reuse

$$\frac{\partial \lambda}{\partial x} = \frac{\partial u}{\partial x} \frac{\partial f}{\partial u} \frac{\partial \lambda}{\partial f}$$

Example



$$\frac{\partial f}{\partial x} = 1 \cdot 1 \cdot 1 = 1$$

$$\frac{\partial f}{\partial y} = 1 \cdot 1 \cdot 1 = 1$$

$$\frac{\partial f}{\partial z} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial z} = (x+y) \cdot 1 = 5$$

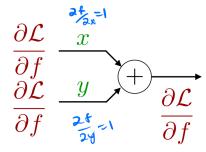
$$x' = 2 \cdot 1 = 12$$

$$y' = 3 \cdot 1 = 22$$

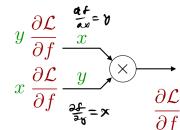
$$z' = 1 \cdot 5 = 3$$

Gradient Gates

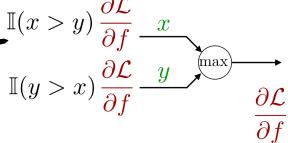
Add: Distribute gradient



Multiply: Switches gradient



Max: Routes gradient



Convergence: Lower Total derivatives, sum all

Multivariate Goods + Chain Rule

Jacobian (vector-vector) $\nabla_x y = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_m} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix}$

$y \in \mathbb{R}^n \quad x \in \mathbb{R}^m$

$j = (\nabla_x y)^T$

$W \in \mathbb{R}^{h \times m} \quad x \in \mathbb{R}^m \quad y \in \mathbb{R}^h$

$$\nabla_x x^T y = y$$

$$\nabla_x Wx = W^T$$

$\nabla_w Wx$ ought to be X^T

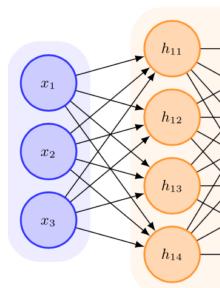
Correct (mostly), only due to sparse matrix

Back Prop NN-Layer

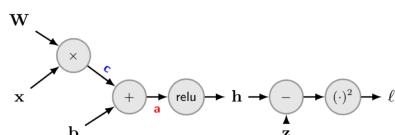
$$y = Wx$$

$$\frac{\partial R}{\partial x} = W^T \frac{\partial E}{\partial y}$$

$$\frac{\partial E}{\partial W} = \frac{\partial E}{\partial y} x^T$$



Here, $h = \text{ReLU}(Wx + b)$, with $h \in \mathbb{R}^h$, $x \in \mathbb{R}^m$, $W \in \mathbb{R}^{h \times m}$, and $b \in \mathbb{R}^h$. While many output cost functions might be used, let's consider a simple squared-loss output $\ell = (h - z)^2$ where z is some target value.

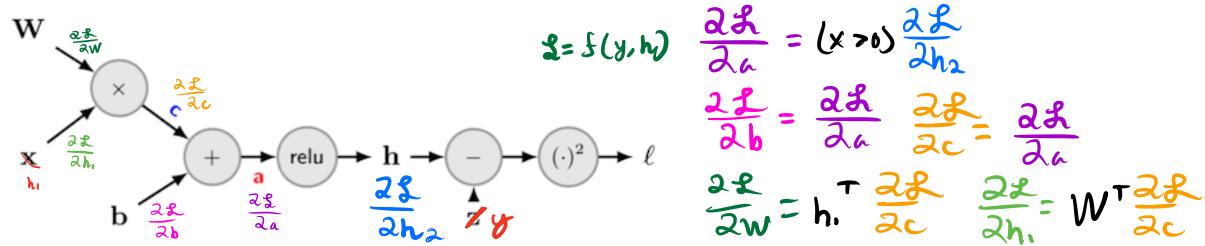


$$\nabla_h \ell = 2(h - z), \quad \frac{\partial \ell}{\partial a} = I(a > 0) \circ \frac{\partial \ell}{\partial h},$$

$$\frac{\partial \ell}{\partial c} = \frac{\partial \ell}{\partial a} \frac{\partial a}{\partial c}, \quad \frac{\partial \ell}{\partial x} = \frac{\partial \ell}{\partial a} \frac{\partial a}{\partial x} = W^T \frac{\partial \ell}{\partial h},$$

$$\frac{\partial \ell}{\partial W} = \frac{\partial \ell}{\partial h} x^T$$

Backprop Basic NN Layer:



Initializations:

- Too small - decline to zero , Too large - explode to infinity

Xavier Init:

$$w_i \sim \mathcal{N}(0, \frac{2}{n_{in} n_{out}})$$

Derivation:

$$\text{Var}(wh) = \sum \text{Var}(h_i) + \sum \text{Var}(w_i) + \text{Var}(w) \text{Var}(h)$$

$$\text{Var}(h_i) = \text{Var}(h_{i-1}) \cdot \sum_{j=1}^{n_h} \text{Var}(w_{ij}) \quad h_{ij} = \sum w_{ij} \cdot h_{i-1,j}$$

since $\text{Var}(h_i) = \text{Var}(h_{i-1})$, need $\sum \text{Var}(w) = 1$, Thus

$$\text{Var}(w) = \frac{1}{n_{in}}$$

Same for backprop and norm

Batch Norm:

For inputs randomly sample to ensure similar distribution / mini batch.

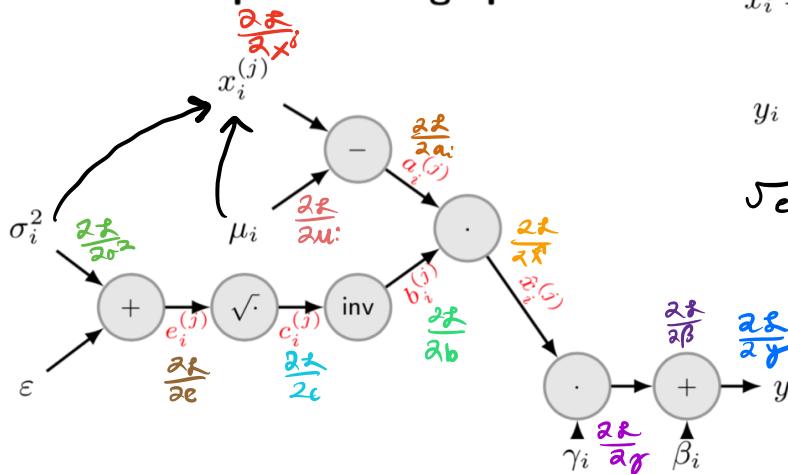
For hidden layers, we use batch norm, solve for internal covariate shift

$$x_i^b = \frac{x_i - \bar{x}_i}{\sqrt{\sigma_i^2 + \epsilon}} \quad \left\{ \begin{array}{l} \bar{x}_i = \frac{1}{m} \sum x_j^i \\ \sigma_i^2 = \frac{1}{m} \sum (x_i - \bar{x}_i)^2 \end{array} \right. \quad \begin{array}{l} \text{across all} \\ \text{samples } (m) \\ \text{in mini batch} \end{array}$$

$$\text{Scale + shift parameters} \quad y_i = \gamma \hat{x}_i + \beta_i$$

Batch Norm Backwards Pass

Batch normalization computational graph



$$\hat{x}_i = \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

$$y_i = \gamma_i \hat{x}_i + \beta_i$$

$$\sqrt{\epsilon} \quad \gamma_i e^{-\frac{1}{2}}$$

$$\frac{\partial \mathcal{L}}{\partial \beta} = \frac{\partial y}{\partial \beta} \quad \frac{\partial \mathcal{L}}{\partial y_i} = 1 \cdot \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial y_i}$$

$$\frac{\partial \mathcal{L}}{\partial \gamma} = \frac{\partial y}{\partial \gamma} \cdot \frac{\partial \mathcal{L}}{\partial y_i} = \sum_{i=1}^m x_i \cdot \frac{\partial \mathcal{L}}{\partial y_i}$$

$$\frac{\partial \mathcal{L}}{\partial x} = \frac{\partial y}{\partial x} \cdot \frac{\partial \mathcal{L}}{\partial y_i} = \frac{\partial \mathcal{L}}{\partial y_i} \cdot \gamma$$

$$\frac{\partial \mathcal{L}}{\partial \alpha} = \frac{\partial x}{\partial \alpha} \cdot \frac{\partial \mathcal{L}}{\partial x} = \frac{1}{\sqrt{\sigma^2 + \epsilon}} \frac{\partial \mathcal{L}}{\partial x}$$

$$\frac{\partial \mathcal{L}}{\partial \mu} = \frac{\partial x}{\partial \mu} \cdot \frac{\partial \mathcal{L}}{\partial x} = -\frac{1}{\sqrt{\sigma^2 + \epsilon}} \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial x} \quad \text{w contributes to all } j's$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial x}{\partial b} \cdot \frac{\partial \mathcal{L}}{\partial x} = (x - \mu) \frac{\partial \mathcal{L}}{\partial x}$$

$$\frac{\partial \mathcal{L}}{\partial c} = \frac{\partial x}{\partial c} \cdot \frac{\partial \mathcal{L}}{\partial x} = -\frac{1}{\sigma^2 + \epsilon} \cdot \frac{\partial \mathcal{L}}{\partial x}$$

$$\frac{\partial \mathcal{L}}{\partial \epsilon} = \frac{\partial x}{\partial \epsilon} \cdot \frac{\partial \mathcal{L}}{\partial x} = \frac{1}{2} \frac{1}{\sqrt{\sigma^2 + \epsilon}} \cdot \frac{\partial \mathcal{L}}{\partial x}$$

$$\frac{\partial \mathcal{L}}{\partial \sigma^2} = \frac{\partial x}{\partial \sigma^2} \cdot \frac{\partial \mathcal{L}}{\partial x} = 1 \cdot \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial x_i}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial x_i} &= \frac{\partial x}{\partial x_i} \cdot \frac{\partial \mathcal{L}}{\partial x} + \frac{\partial x^2}{\partial x_i} \cdot \frac{\partial \mathcal{L}}{\partial x^2} + \frac{\partial \mu}{\partial x_i} \cdot \frac{\partial \mathcal{L}}{\partial \mu} \\ &= 1 \cdot \frac{1}{\sqrt{\sigma^2 + \epsilon}} \frac{\partial \mathcal{L}}{\partial x} + \frac{2}{m} (x_i - \mu) \frac{\partial \mathcal{L}}{\partial x^2} + \frac{1}{m} \frac{\partial \mathcal{L}}{\partial \mu} \end{aligned}$$

Lec 9 :

2/2

Regularization - any modification to reduce generalization error, but not training error.

- Increase estimator bias but decrease variance
- Prevent overfitting

Parameter Norm -

$$J(\theta) + \alpha \Omega(\theta)$$

$\|w\|_2$ L2 Norm

$\|w - b\|_1$ Manhattan

$\|w_1, w_2\|_1$ similar features

$\|w\|_1$, L1 - sparse solutions

$\|h\|_1$, Sparse representation

Multi-task learning: Composed NN's

Transfer learning: Add a head layer(s) with little training on existing model for similar tasks

Ensemble Methods: Train many models then avg

Bayesian - k of M samples, k models, avg (M total sample $M \gg N$)

Dropout: Mask nodes with fixed prob. Approximates ^{model} bayesian in single

Sample a binary mask for all units then do FW, BW pass.

New mask for every SGD step. 2^N configs

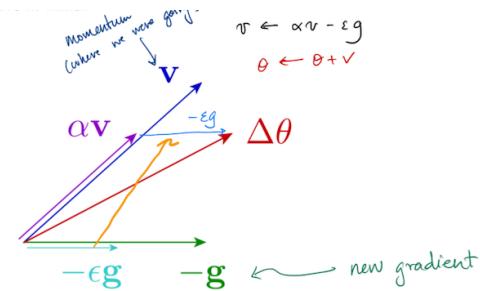
- Approximates bayesian in one model
- Regularizes each hidden unit to work in many contexts
 - May encode redundancy
- Scale by $1/p$ in training so output is appropriately scaled.

Optimization (GD variations):

Momentum

$$\text{Momentum : } \begin{aligned} V &\leftarrow \alpha V - \epsilon g & \text{learn rate} \\ \theta &\leftarrow \theta + V \end{aligned}$$

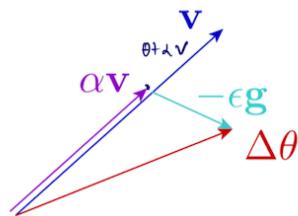
Pros: Avg's out back and forth, can pull us out of local optima



Nesterov Momentum

$$\begin{aligned} V &\leftarrow \alpha V + \epsilon \nabla_{\theta} J(\theta + \alpha V) \\ \theta &\leftarrow \theta + V \end{aligned}$$

Pros: Gradient computed after momentum step so more accurate.



Adagrad!

Scales grad down over time ... dows

$$\begin{aligned} a &\leftarrow a + g \odot g \\ \theta &\leftarrow \theta - \frac{\epsilon}{\sqrt{a}} \odot g \end{aligned}$$

RMS Prop:

Weighted version of adagrad

$$\begin{aligned} a &\leftarrow \alpha a + (1-\beta) g \odot g & \text{component wise} \\ \theta &\leftarrow \theta - \frac{\epsilon}{\sqrt{a} + \nu} \odot g \end{aligned}$$

Adam:

$$\text{1st moment momentum } v = \beta_1 v + (1-\beta_1) g$$

$$\text{2nd moment grad norm } a = \beta_2 a + (1-\beta_2) g \odot g$$

$$\beta^{\text{bias correction}} \quad \tilde{v} = \frac{1}{1-\beta_1} v \quad \tilde{a} = \frac{1}{1-\beta_2} a$$

$$\text{step!} \quad \theta \leftarrow \theta - \frac{\epsilon}{\sqrt{\tilde{a}} + \nu} \odot \tilde{v}$$

Convolution

In practice, it's cross corr, or just diag (not full filter)

- typically depth matters

Parameters: $(h \times w \times d) \times \#r$

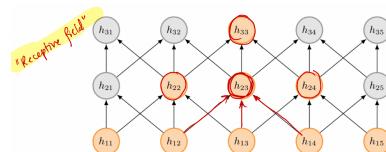
Spatiality:

CNN: $w \times h \times d$ weights / neuron

FC: $w \cdot h \cdot d$ weights / neuron

Receptive Field:

of inputs a neuron 'sees'



Padding: Add zeros to maintain output size

Stride: Skip on diag of filter convolution

$$\left(\frac{w - w_f + 2\text{pad}}{\text{stride}} + 1, \frac{h - h_f + 2\text{pad}}{\text{stride}} + 1 \right)$$

Pooling layer: Down sampling with specified operation stride
No parameters, deal with parallel dimensions

$$\left(\frac{w - w_p}{\text{stride}} + 1, \frac{h - h_p}{\text{stride}} + 1 \right)$$

