## Convolution:

$$S(t) = \int x(a)\, w(t-a)\, da$$

$\uparrow$ (feature map) $= (x * w)(t)$

"Input" "Kernel"

| | |
|---|---|
| Commutative Property | $(A*B)(t)$ $\equiv$ $(B*A)(t)$ |

$$S(i,j) = (K * I)(i,j) = \sum_m \sum_n I(i+m, j+n)\, K(m,n)$$

## 2D Discrete Convolution:

$$S(i,j) = (I*K)(i,j) = \sum_m \sum_n I(m,n)\, K(i-m, j-n)$$

$$\equiv (K*I)(i,j) = \sum_m \sum_n I(i-m, j-n)\, K(m,n)$$

$= Conv$ (Output)

$$\begin{array}{|c|c|c|} \hline \times_{11} & \times_{12} & \times_{13} \\ \hline \times_{21} & \times_{22} & \times_{23} \\ \hline \times_{31} & \times_{32} & \times_{33} \\ \hline \end{array} \quad \begin{array}{|c|c|} \hline F_{11} & F_{12} \\ \hline F_{21} & F_{22} \\ \hline \end{array}$$

$\begin{array}{|c|c|} \hline O_{11} & O_{12} \\ \hline O_{21} & O_{22} \\ \hline \end{array}$

where $\begin{cases} X = \text{Input} \\ F = \text{Filter} \end{cases}$

& $O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22}$

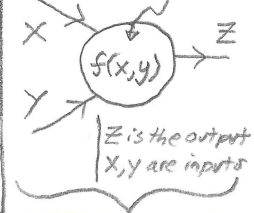$\boxed{\begin{array}{c} \text{Input Filter Output} \\ (\text{Convolution of } X * F) \end{array}}$

## CHAIN RULE OF DIFFERENTIATION

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \times \frac{\partial g}{\partial x}$$

## CONVOLUTION FILTER UPDATE

$$\boxed{F_{updated} = F - \alpha \frac{\partial L}{\partial F}}$$

## COMPUTATIONAL GRAPH

$X$, $Y$ are inputs; $Z$ is the output. $f(x,y)$ (gate)

## FORWARD PASS

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial x}$$ (L is the Loss fn)

Local Gradient $\nabla f$ (gradients) $\frac{\partial L}{\partial z}$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial y}$$

## BACKWARD PASS

$\frac{\partial L}{\partial z} = $ "Loss gradient from previous layer"

$\frac{\partial z}{\partial y} = $ "Local Gradients"

$\frac{\partial L}{\partial y} = $ "Gradient to update the filter with"
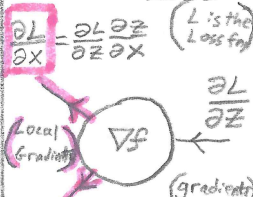
---

## full Convolution ← ZERO-Padding

$$\frac{\partial L}{\partial F} = Convolution\left(X, \frac{\partial L}{\partial O}\right)$$ (Loss gradient)

$$\frac{\partial L}{\partial X} = Full\,Convolution\left(\begin{array}{c}F\,Rotated \\ by\,180°\end{array}, \frac{\partial L}{\partial O}\right)$$ (Input Image)

**Typical Layout:** Input → Convolution → ReLU → Pooling → FC

## Dropout (with probability $P$ that unit is present during training)

(dropout) → (for training)

$x' = \frac{x}{|x|}$

$* $ Bernoulli Distribution
$$r_j^l \sim Bernoulli(p) = \begin{cases} 1 & \text{for KEEP} \\ 0 & \text{for drop} \end{cases}$$

$\tilde{y}^{(l)} = r^{(l)} * y^{(l)}$ (* indicates dropout)
$z_i^{(l+1)} = w_i^{(l+1)} \tilde{y}^{(l)} + b_i^{(l+1)}$

$*$ When testing, all units remain & weights $\omega$ are multiplied by prob. $p$

$P \sim \begin{cases} 0.5 \text{ for hidden layer} \\ 0.8 \text{ for input layer} \end{cases}$

$*$ Element-Wise Multiplication $A \odot B$ or $A \otimes B$

$P \to PW$ (Training) (Testing)

$\Rightarrow$ **PREVENTS OVERFITTING**

## Dropout vs Ridge Regression (L2)

Standard Network → Dropout Network

## Convolution Layer:

filter map: $n \times n \times q$ (Kernel size) (filter count, represents different colors)

feature map: $m - n + 1$

$m \times m \times r$ (size) (channel)

$\Rightarrow$ Input convolved w/ filter results in a feature map

## BACKPROPAGATION WITH CONVOLUTION LAYERS

$$\delta^{(l)} = W^{(l)t} \delta^{(l+1)} \cdot f'(z^{(l)})$$ (transpose)
(for dense connections)
($l$ = convolutional layer)

Cost Function: $J(w, b, x, y)$ (training data)

Error Term: $\delta^{(l+1)}$ (of the (l-1)-th layer)

$*$ If conv layer is sub-sampling layer too (i.e., mean or max pooling)

$$S_k^{(l)} = upsample\left(W_k^{(l)t} \delta_k^{(l+1)}\right) \cdot f'(z_k^{(l)})$$
(where $K$ is the filter index)

$\begin{cases} P = \text{Zero-Padding} \\ S = \text{stride} \end{cases}$

## Ridge Regression (L2 Regularization):

$$minimize \left\{ |y - X\omega|^2 + \lambda |\omega|^2 \right\}$$ (target)

$*$ Smooths out behaviour of results

(forces $\omega$ norms to be small)

for dropout the result to minimize is
$$\left\{ |y - pX\omega|^2 + P(1-P)|\Gamma\omega|^2 \right\}$$
where $\Gamma = diag(X^tX)$
$\Rightarrow$ SIMILAR RESULTS

## Partial Derivatives of a matrix $O$ (output) wrt matrix $F$ (filter)

$$\frac{\partial L}{\partial F_i} = \sum_{k=1}^{M} \frac{\partial L}{\partial O_k} \frac{\partial O_k}{\partial F_i}$$
(Previous Layer) (Local grad)

(via chain rule)

**OUTPUT SIZE** $= (W - F + 2P)/S + 1$
maintain spatial size by $\begin{cases} P = \frac{1}{2}(F-1) \\ S = 1 \end{cases}$

## Bernoulli Random Variable $(X)$

$X = \begin{cases} 1: \text{Success} \\ 0: \text{Failure} \end{cases}$ & PMF follows:

$S_x = \{0, 1\}$

$P_x(1) = P, \; P_x(0) = 1-p$

$\mu_x = E[X] = p, \; \sigma_x^2 = Var[X] = p(1-p)$

## KL Divergence (Kullback-Leibler)

$\sum_{j=1}^{S_2} KL(\rho, \hat{\rho}_j)$ $\begin{cases} \text{Divergence btwn a} \\ \text{Bernoulli RV with} \\ \text{means } p \text{ \& } \hat{\rho}_j \end{cases}$

$(S_2 = \text{\# of hidden nodes})$ (if similar, $\approx \emptyset$)

---

\* Note that the "sparsity" of convolutional layer is due to the Parameter sharing of weights

---

\* Directed Acyclic Graph (DAG) Networks

\* ResNet is a type of DAG with a residual connection (shortcut) that bypasses main network layers



Input
Conv / Batch Norm / ReLU / Pool
(Skip layer connection)
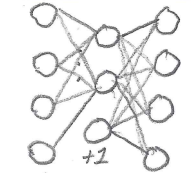(+) (addition)
ReLU
Output

---

\* Good way to normalize data is to subtract the mean & divide by the standard deviation along each feature (higher accuracy)

$X' \rightarrow \dfrac{X - \mu_x}{\sigma_x}$

\* DONE PER FEATURE (BatchNorm)

---

## Autoencoders (for data compression, denoising, feature maps, classification, sparse representation)



Hidden Representation
UNSUPERVISED TRAINING

\* Aim is for output to equal input
\* Consider hidden node: $a_j^2$ (activation of hidden unit $j$)
\* Given the input $X$: $a_j^2(x^{(i)})$ is the current activation of the hidden units $j$ after receiving input $X$ (corresponding to example) in the training set

FEATURE EXTRACTION
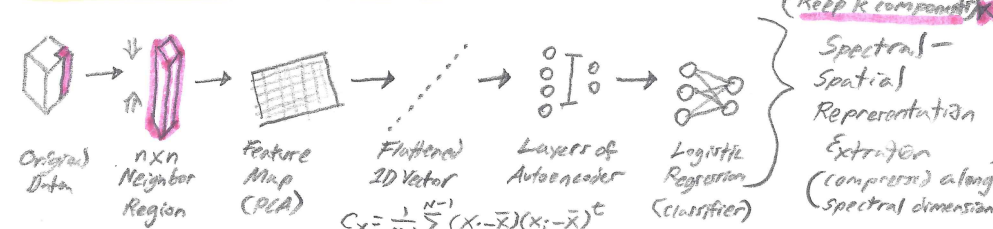
## Batch Normalization (done in minibatch of # of vectors)

(mean)   (vectors)

$y = \left( \dfrac{X - E[X]}{\sqrt{Var[X] + Constant}} \right) \gamma + \beta$

(to avoid numerical instability)

Faster Convergence
Regularization
Decreased Importance of weight initialization

\* Standard Normalization (most common): $\gamma$ is unity & $\beta = \emptyset$

## Principal Component Analysis (PCA): "Converting input vector X into feature vector Y"



Original Data → n×n Neighbor Region → Feature Map (PCA) → Flattened 1D Vector → Layers of Autoencoder → Logistic Regression (classifier)

keep k components
Spectral-Spatial Representation Extraction (compressed along spectral dimension)

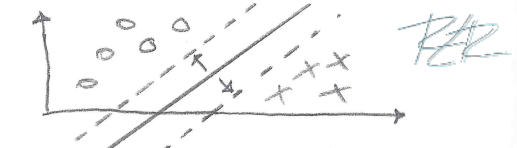$C_x = \dfrac{1}{N-1} \sum^{N-1} (X_i - \bar{X})(X_i - \bar{X})^t$

## DKLT Theory: Transformation $T$ to diagonalize the covariance matrix of a random signal in the discrete time domain for signal compression
$\Rightarrow$ PCA Approximation through uncorrelated elements of $Y$ (keeping M components)

## SVM-Kernel Functions (eg., RBFNs)

Kernel function (usually non linear) $\}$ $f\left( \underbrace{|x - \omega|^2}_{\text{data} \;\; \text{params}} \right)$ (distance$^2$)

Example: Gaussian Kernel Function:

$f \sim e^{-|x - \omega|^2 / 2\sigma^2}$

$\phi \leftarrow f$ (for kernel functions)

$\Rightarrow \phi$ returns inner product of two points in a feature space

Polynomial Kernel: (Image Processing)
$K(x_i, x_j) = (x_i \cdot x_j + 1)^d$

Gaussian Kernel: (NO PRIOR KNOWLEDGE)
$K(x, y) = \exp\left( \dfrac{-||x-y||^2}{2\sigma^2} \right)$

Gaussian Radial Basis Function: (NO PRIOR KNOWLEDGE)
$K(x_i, x_j) = \exp(-\gamma ||x_i - x_j||^2)$

Kernel Functions (SVMs) take input vectors in original space & returns their dot product in the feature space

$K(x, x') = \langle \phi(x), \phi(x') \rangle$

\* Mapping from input space to a dual space ("Gram Matrix")

## Sparsity Parameter $\rho$ (constraint) (typically $\rho \approx 0.05$)

\* ENFORCE APPROXIMATION CONSTRAINT

$\hat{\rho}_j = \rho$ (requires an extra penalty term in optimization)

where $\hat{\rho}_j = \dfrac{1}{m} \sum_{i=1}^{m} a_j^2(x^{(i)})$

Average activation for hidden unit $j$ over the entire training data set (containing $m$ examples)

---

\* Support Vector Machines (SVMs): for linearly seperable decision boundaries defined as a border with as much space between vector classes as possible



(if $K = N \leftarrow$ NO LOSS OF INFORMATION)
Eigenvalues with large variance are the most important components to keep (i.e., keep $\lambda$ if large)

\* for real signals: $\rho \approx 0.9$ (speech modeling)

$\lambda_k = \dfrac{1 - \rho^2}{1 - 2\rho \cos(w_k) + \rho^2}$

## Discrete Cosine Transform (DCT)
is a good approximation of KLT

## Toeplitz Matrix (or pL2)
(forgets anything except 1 element)

$C_x = \begin{bmatrix} 1 & \rho & \rho^2 \\ \rho & 1 & \rho \\ \rho^2 & \rho & 1 \end{bmatrix}$ $\begin{cases} \text{Covariance Matrix for} \\ \text{a zero-mean, first-} \\ \text{order Markov sequence} \\ \text{(of } N \text{ elements)} \end{cases}$

$X_n \leftarrow \rho X_{n-1} + \epsilon$