

INT305 note

(Machine Learning)

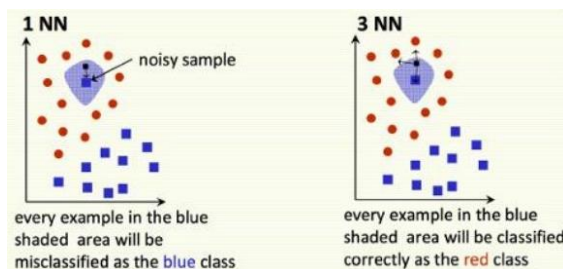
1 Introduction

1.1 Supervised learning (much of this course)

Task	Inputs	Labels
object recognition	image	object category
image captioning	image	caption
document classification	text	document category
speech-to-text	audio waveform	text
\vdots	\vdots	\vdots

1.1.1 KNN

- Nearest neighbours **sensitive to noise or mis-labelled data** ("class noise").
- Smooth by having k nearest neighbours vote.

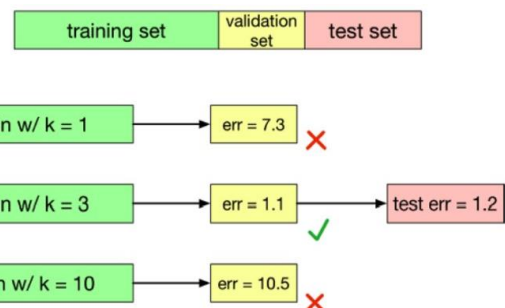
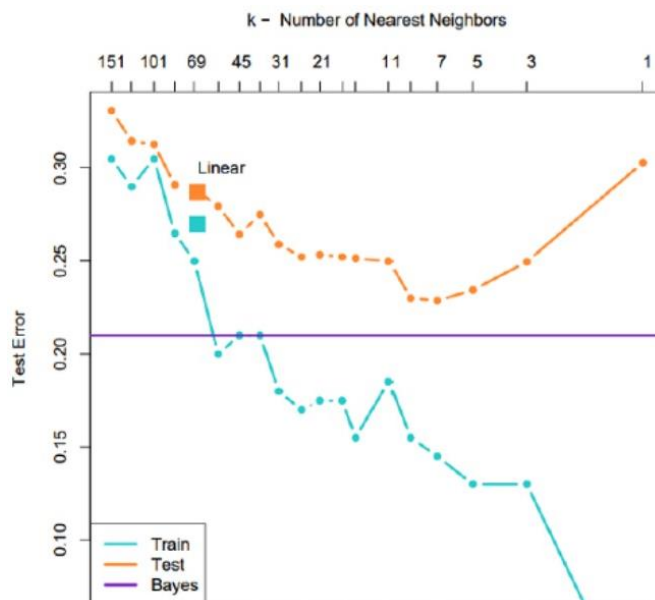


Algorithm (kNN):

1. Find k examples $\{\mathbf{x}^{(i)}, t^{(i)}\}$ closest to the test instance \mathbf{x}
2. Classification output is majority class

$$y = \arg \max_{t^{(z)}} \sum_{i=1}^k \mathbb{I}(t^{(z)} = t^{(i)})$$

- Balancing hyperparameter k
 - Optimal choice of k depends on number of data points n .
 - Nice theoretical properties if $k \rightarrow \infty$ and $k/n \rightarrow 0$.
 - Rule of thumb: choose $k < \sqrt{n}$.
 - We can choose k using validation set.



2 Linear Methods for Regression, Optimization

Linear regression exemplifies recurring themes of this course:

- Choose a **model** and a **loss function**
- Formulate an **optimization problem**
- Solve the minimization problem using one of two strategies

➤ **Direct solution** (set derivatives to zero)

➤ **Gradient descent**

- **Vectorize** the algorithm, i.e. represents in terms of linear algebra
- Make a linear model more powerful using features
- Improve the generalization by adding a **regularizer**

2.1 Supervised Learning Setup

In supervised learning:

- There is input $\mathbf{x} \in \mathcal{X}$, typically a vector of features (or covariates)
- There is target $t \in \mathcal{T}$, (also called response, outcome, output, class)
- Objective is to learn a function $f: \mathcal{X} \rightarrow \mathcal{T}$ such that $t \approx y = f(\mathbf{x})$ based on some data $\mathcal{D} = \{(\mathbf{x}^{(i)}, t^{(i)}) \text{ for } i = 1, 2, \dots, N\}$

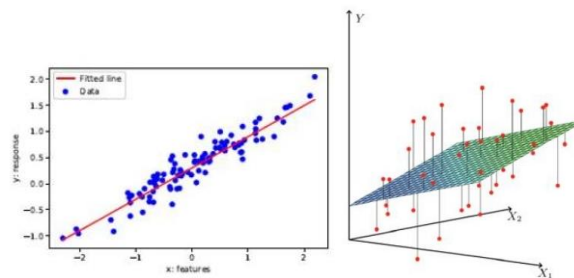
2.2 Linear Regression

2.2.1 Linear Regression Model

Model: In linear regression, we use a linear function of the features $\mathbf{x} = x_1, \dots, x_D \in \mathbb{R}^D$ to make predictions y of the target value $t \in \mathbb{R}$:

$$y = f(\mathbf{x}) = \sum_j w_j x_j + b$$

- y is the **prediction**
- \mathbf{w} is the **weights**
- b is the **bias** (or intercept)
- \mathbf{w} and b together are the **parameters**
- We hope that our prediction is close to the target: $y \approx t$.



- If we have only 1 feature: $y = wx + b$ where $w, x, b \in \mathbb{R}$.
- y is linear in x .
- If we have only D features: $y = \mathbf{w}^T \mathbf{x} + b$ where $\mathbf{w}, \mathbf{x} \in \mathbb{R}^D, b \in \mathbb{R}$
- y is linear in \mathbf{x} .

2.2.2 Linear Regression workflow

We have a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, t^{(i)}) \text{ for } i = 1, 2, \dots, N\}$:

- $\mathbf{x}^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_D^{(i)})^T \in \mathbb{R}^D$ are the inputs (e.g. age, height)
- $t^{(i)} \in \mathbb{R}$ is the target or response (e.g. income)
- Predict $t^{(i)}$ with a linear function of $\mathbf{x}^{(i)}$:
 - $t^{(i)} \approx y^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + b$
 - Different \mathbf{w}, b define different lines.
 - We want the "best" line \mathbf{w}, b .

2.2.3 Linear Regression Loss Function

• A loss function $\mathcal{L}(y, t)$ defines how bad it is if, for some example \mathbf{x} , the algorithm predicts y , but the target is actually t .

• **Squared error loss function:**

$$\mathcal{L}(y, t) = \frac{1}{2} (y - t)^2$$

- $y - t$ is the residual, and we want to make this small in magnitude.
- The $1/2$ factor is just to make the calculations convenient.

- **Cost function:** loss function averaged over all training examples.

$$\mathcal{J}(\mathbf{w}, b) = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - t^{(i)})^2 = \frac{1}{2N} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}^{(i)} + b - t^{(i)})^2$$

2.2.3 Linear Regression Vectorization

But if we expand $y^{(i)}$, it will get messy:

$$\frac{1}{2N} \sum_{i=1}^N \left(\sum_{j=1}^D (w_j x_j^{(i)}) + b - t^{(i)} \right)^2$$

Vectorize algorithms by expressing them in terms of vectors and matrices:

$$\mathbf{w} = (w_1, \dots, w_D)^T \quad \mathbf{x} = (x_1, \dots, x_D)^T$$

$$y = \mathbf{w}^T \mathbf{x} + b$$

Python code:

```

y = b
for j in range(M):
    y += w[j] * x[j]
    
```

$$=$$

```

y = np.dot(w, x) + b
    
```

Organize all the training examples into a **design matrix X** with one row per training example, and all the targets into the **target vector t**:

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}^{(1)T} \\ \mathbf{x}^{(2)T} \\ \mathbf{x}^{(3)T} \end{pmatrix} = \begin{pmatrix} 8 & 0 & 3 & 0 \\ 6 & -1 & 5 & 3 \\ 2 & 5 & -2 & 8 \end{pmatrix}$$

one feature across
all training examples

one training
example (vector)

Computing the **predictions** for the whole dataset:

$$\mathbf{X}\mathbf{w} + b\mathbf{1} = \begin{pmatrix} \mathbf{w}^T \mathbf{x}^{(1)} + b \\ \vdots \\ \mathbf{w}^T \mathbf{x}^{(N)} + b \end{pmatrix} = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{pmatrix} = \mathbf{y}$$

Computing the **squared error cost** across the whole dataset:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + b\mathbf{1}$$

$$\mathcal{J} = \frac{1}{2N} \|\mathbf{y} - \mathbf{t}\|^2$$

We can also add a column of 1's to design matrix, combine the bias and the weights, and conveniently write:

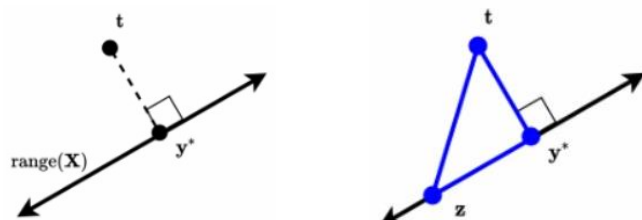
$$\mathbf{X} = \begin{bmatrix} 1 & [\mathbf{x}^{(1)}]^T \\ 1 & [\mathbf{x}^{(2)}]^T \\ \vdots & \vdots \end{bmatrix} \in \mathbb{R}^{N \times (D+1)} \quad \text{and} \quad \mathbf{w} = \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \end{bmatrix} \in \mathbb{R}^{D+1}$$

Then, our predictions reduce to $\mathbf{y} = \mathbf{X}\mathbf{w}$.

2.3 Direct Solution •

2.3.1 Linear Algebra •

- We seek \mathbf{w} to minimize $\|\mathbf{X}\mathbf{w} - \mathbf{t}\|^2$, or equivalently $\|\mathbf{X}\mathbf{w} - \mathbf{t}\|$
- $\text{range } \mathbf{X} = \{\mathbf{X}\mathbf{w} | \mathbf{w} \in \mathbb{R}^D\}$ is a D -dimensional subspace of \mathbb{R}^N
- Recall that the closest point $\mathbf{y}^* = \mathbf{X}\mathbf{w}^*$ in subspace $\text{range}(\mathbf{X})$ of \mathbb{R}^N to arbitrary point $\mathbf{t} \in \mathbb{R}^N$ is found by orthogonal projection.
- We have $(\mathbf{y}^* - \mathbf{t}) \perp \mathbf{X}\mathbf{w}, \forall \mathbf{w} \in \mathbb{R}^D$
- \mathbf{y}^* is the closest point to \mathbf{t}



2.3.2 Calculus •

- **Partial derivative**: derivative of a multivariate function with respect to one of its arguments.

$$\frac{\partial}{\partial x_1} f(x_1, x_2) = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2) - f(x_1, x_2)}{h}$$

- To compute, take the single variable derivative, pretending the other arguments are constant.
- Example: partial derivatives of the prediction y .

$$\begin{aligned} \frac{\partial y}{\partial w_j} &= \frac{\partial}{\partial w_j} \left[\sum_{j'} w_{j'} x_{j'} + b \right] & \frac{\partial y}{\partial b} &= \frac{\partial}{\partial b} \left[\sum_{j'} w_{j'} x_{j'} + b \right] \\ &= x_j & &= 1 \end{aligned}$$

- For **loss derivatives**, apply the chain rule:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_j} &= \frac{d\mathcal{L}}{dy} \frac{\partial y}{\partial w_j} & \frac{\partial \mathcal{L}}{\partial b} &= \frac{d\mathcal{L}}{dy} \frac{\partial y}{\partial b} \\ &= \frac{d}{dy} \left[\frac{1}{2} (y - t)^2 \right] \cdot x_j & &= y - t \\ &= (y - t) x_j & & \end{aligned}$$

- For cost derivatives, use linearity and average over data points:

$$\frac{\partial \mathcal{J}}{\partial w_j} = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)} \quad \frac{\partial \mathcal{J}}{\partial b} = \frac{1}{N} \sum_{i=1}^N y^{(i)} - t^{(i)}$$

- Minimum must occur at a point where partial derivative are zero.

$$\frac{\partial \mathcal{J}}{\partial w_j} = 0 \quad (\forall j), \quad \frac{\partial \mathcal{J}}{\partial b} = 0.$$

(If $\partial \mathcal{J} / \partial w_j \neq 0$, you could reduce the cost by changing w_j)

- We call the vector of partial derivatives the **gradient**.

- Thus, the "gradient of $f: \mathbb{R}^D \rightarrow \mathbb{R}$ ", denoted $\nabla f(\mathbf{w})$, is:

$$\left(\frac{\partial}{\partial w_1} f(\mathbf{w}), \dots, \frac{\partial}{\partial w_D} f(\mathbf{w}) \right)^\top$$

- The gradient points in the direction of the greatest rate of increase.
- Analogue of second derivative (the "Hessian" matrix):

$$\nabla^2 f(\mathbf{w}) \in \mathbb{R}^{D \times D} \text{ is a matrix with } [\nabla^2 f(\mathbf{w})]_{ij} = \frac{\partial^2}{\partial w_i \partial w_j} f(\mathbf{w})$$

- We seek \mathbf{w} to minimize $\mathcal{J}(\mathbf{w}) = ||\mathbf{X}\mathbf{w} - \mathbf{t}||^2 / 2$
- Taking the gradient with respect to \mathbf{w} we get:

$$\nabla_{\mathbf{w}} \mathcal{J}(\mathbf{w}) = \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{t} = 0$$

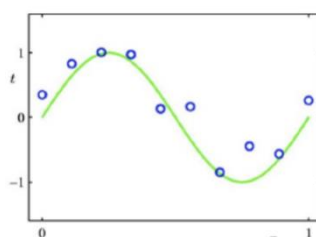
$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

- **Linear regression** is one of **only** a handful of models in this course that **permit direct solution**.

2.4 Polynomial Feature Mapping

2.4.1 Introduction

The relation between the input and output may not be linear. But we can still use linear regression by mapping the input features to another space using **feature mapping** (or **basis expansion**). $\varphi(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and treat the mapped feature (in \mathbb{R}^d) as the input of a linear regression procedure.



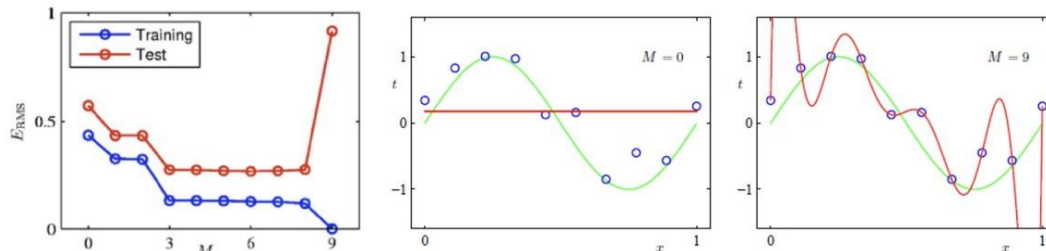
Find the data using a degree-M polynomial function of the form:

$$y = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{i=0}^M w_i x^i$$

- Here the feature mapping is $\varphi(x) = [1, x, x^2, \dots, x^M]^T$.
- We can use linear regression to find \mathbf{w} , since $y = \varphi(x)^T \mathbf{w}$ is linear with w_0, w_1, \dots, w_M

2.4.2 Model Complexity and Generalization

- Underfitting (M=0): model is too simple – does not fit the data.
- Overfitting (M=9): model is too complex – fits perfectly.



2.4.3 L² Regularization

Regularizer: a function that quantifies how much we prefer one hypothesis VS another.

- We can encourage the weights to be small by choosing as our regularizer the **L² penalty**.

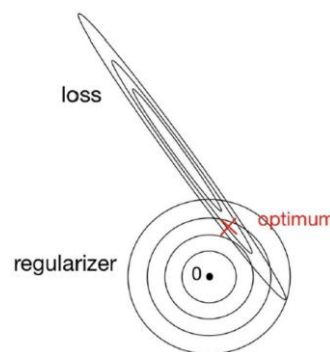
$$\mathcal{R}(\mathbf{w}) = \frac{1}{2} \sum_j w_j^2$$

(To be precise, the L² norm is Euclidean distance, we're regularizing the squared L² norm)

- The **regularized cost function** makes a tradeoff between fit to the data and the norm of the weights.

$$J_{reg}(\mathbf{w}) = J(\mathbf{w}) + \lambda \mathcal{R}(\mathbf{w}) = J(\mathbf{w}) + \frac{\lambda}{2} \sum_j w_j^2$$

- If you fit training data poorly, J is large. If your optimal weights have high values, \mathcal{R} is large.
- Large λ penalize wight values more.
- Like M , λ is a hyperparameter we can tune with a validation set.



2.4.4 L² Regularized Least Squares: Ridge regression

For the least squares problem, we have $J(\mathbf{w}) = \frac{1}{2N} \|\mathbf{X}\mathbf{w} - \mathbf{t}\|^2$

- When $\lambda > 0$ (with regularization), regularized cost gives:

$$\begin{aligned} \mathbf{w}_\lambda^{\text{Ridge}} &= \underset{\mathbf{w}}{\operatorname{argmin}} J_{reg}(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2N} \|\mathbf{X}\mathbf{w} - \mathbf{t}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t} \end{aligned}$$

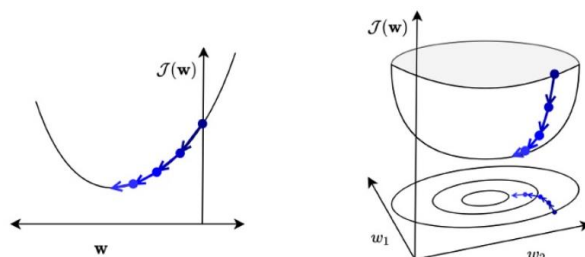
- The case $\lambda = 0$ (no regularization) reduces to least squares solution!

2.5 Gradient Descent

2.5.1 Concepts

- Many times, we do not have a direct solution: Taking derivatives of J w.r.t \mathbf{w} and setting them to 0 doesn't have an explicit solution.

- Gradient descent is an iterative algorithm, which means we apply an update repeatedly until some criterion is met.
- We initialize the weights to something reasonable (e.g., all zeros) and repeatedly adjust them in the direction of steepest descent.



(就是等到斜率为 0 即为最优解)

- Observe:
 - If $\partial J / \partial w_j > 0$, then increasing w_j increases J .
 - If $\partial J / \partial w_j < 0$, then increasing w_j decreases J .
- The following update always decreases the cost function for small enough α (unless $\partial J / \partial w_j = 0$):
- $\alpha > 0$ is a learning rate (or step size). The larger it is, the faster \mathbf{w} changes (but values are typically small).
- This gets its name from the gradient:

$$\nabla_{\mathbf{w}} J = \frac{\partial J}{\partial \mathbf{w}} = \begin{pmatrix} \frac{\partial J}{\partial w_1} \\ \vdots \\ \frac{\partial J}{\partial w_D} \end{pmatrix}$$

- Update rule in **vector form**:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial J}{\partial \mathbf{w}}$$

- And for **linear regression** we have:

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) \mathbf{x}^{(i)}$$

- So gradient descent updates \mathbf{w} in the direction of fastest decrease.
- Observe that once it converges, we get a **critical point**. i.e. $\frac{\partial J}{\partial \mathbf{w}} = \mathbf{0}$

2.5.2 Gradient Descent for Linear Regression

- Why gradient descent, if we can find the optimum directly?
 - gradient descent can be applied to a much broader set of models
 - gradient descent can be easier to implement than direct solutions
 - For regression in high-dimensional space, gradient descent is more efficient than direct solution

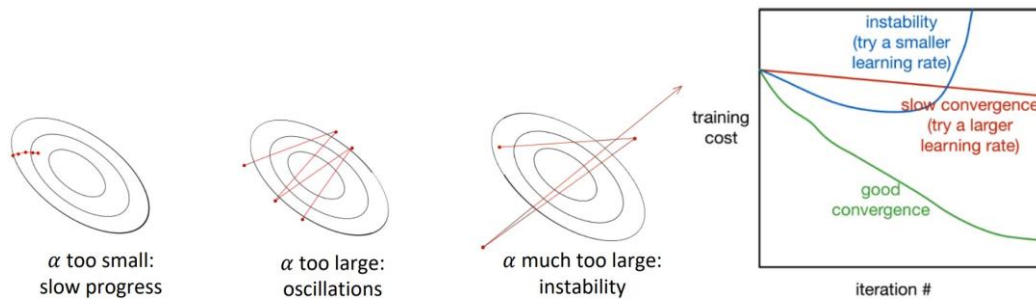
2.5.3 Gradient Descent under the L^2 Regularization

- The gradient descent update to minimize the L^2 regularized cost $J + \lambda \mathcal{R}$ results in **weight decay**:

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \alpha \frac{\partial}{\partial \mathbf{w}} (J + \lambda \mathcal{R}) \\ &= \mathbf{w} - \alpha \left(\frac{\partial J}{\partial \mathbf{w}} + \lambda \frac{\partial \mathcal{R}}{\partial \mathbf{w}} \right) \\ &= \mathbf{w} - \alpha \left(\frac{\partial J}{\partial \mathbf{w}} + \lambda \mathbf{w} \right) \\ &= (1 - \alpha \lambda) \mathbf{w} - \alpha \frac{\partial J}{\partial \mathbf{w}} \end{aligned}$$

2.5.4 Learning Rate (Step Size)

- In gradient descent, the learning rate α is a hyperparameter we need to tune.
- Good values are typically between 0.001 and 0.1.



- To diagnose optimization problems, it's useful to look at **training curves**: plot the training cost as a function of iteration.
- Warning: in general, it's very hard to tell from the training curves whether an optimizer has converged. They can reveal major problems, but they can't guarantee convergence.

2.5.5 Stochastic Gradient Descent

Stochastic gradient descent (SGD): update the parameters based on the gradient for a single training example

1- Choose i uniformly at random

$$2- \theta \leftarrow \theta - \alpha \frac{\partial \mathcal{L}^{(i)}}{\partial \theta}$$

- Cost of each SGD update is independent of N !
- SGD can make significant progress before even seeing all the data!
- Mathematical justification: if you sample a training example uniformly at random, the stochastic gradient is an **unbiased estimate** of the batch gradient:

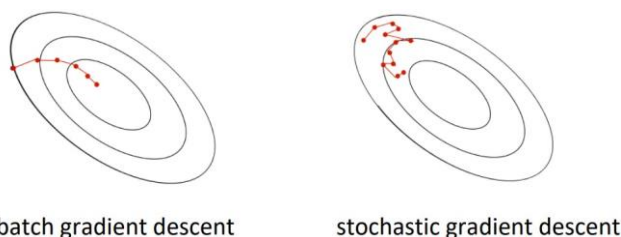
$$\mathbb{E} \left[\frac{\partial \mathcal{L}^{(i)}}{\partial \theta} \right] = \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}^{(i)}}{\partial \theta} = \frac{\partial \mathcal{J}}{\partial \theta}$$

2.5.6 Mini-batch Stochastic Gradient Descent

- **Problems** with using single training example to estimate gradient:
 - Variance in the estimate may be high
 - We can't exploit efficient vectorized operations
- **Compromise approach**:
 - Compute the gradients on a randomly chosen medium-sized set of training examples $\mathcal{M} \subset \{1, \dots, N\}$ called a **mini-batch**.
- Stochastic gradients computed on larger mini-batches have smaller variance.
- The mini-batch size $|\mathcal{M}|$ is a hyperparameter that needs to be set.
 - Too large: requires more compute: e.g., it takes more memory to store the activations, and longer to compute each gradient update
 - Too small: can't exploit vectorization, has high variance
 - A reasonable value might be $|\mathcal{M}| = 100$.

2.5.7 Comparison

- Batch gradient descent moves directly downhill (locally speaking).
- SGD takes steps in a noisy direction, but moves downhill on average.



▲ **Batch Gradient Descent**, 全批量梯度下降, 是最原始的形式, 它是指在每一次迭代时使用所有样本来进行梯度的更新。优点是全局最优解, 易于并行实现; 缺点是当样本数目很大时, 训练过程会很慢。

▲ **Stochastic Gradient Descent**, 随机梯度下降, 是指在每一次迭代时使用一个样本来进行参数的更新。优点是训

练速度快；缺点是准确度下降，并且可能无法收敛或者在最小值附近震荡。

▲ Mini-Batch Gradient Descent, 小批量梯度下降, 是对上述两种方法的一个折中办法。它是指在每一次迭代时使用一部分样本来进行参数的更新。这种方法兼顾了计算速度和准确度。

3 Linear Classifiers, Logistic Regression, Multiclass Classification

3.1 Binary linear classification

- Classification: given a D -dimensional input $\mathbf{x} \in \mathbb{R}^D$ predict a discrete-valued target
- Binary: predict a binary target $t \in \{0,1\}$
 - Training examples with $t = 1$ are called positive examples, and training examples with $t = 0$ are called negative examples.
 - $t \in \{0,1\}$ or $t \in \{-1, +1\}$ is for computational convenience.
- Linear: model prediction y is a linear function of \mathbf{x} , followed by a threshold r :

$$z = \mathbf{w}^T \mathbf{x} + b$$

$$y = \begin{cases} 1 & \text{if } z \geq r \\ 0 & \text{if } z < r \end{cases}$$

- **Eliminating the threshold:** We can assume without loss of generality (WLOG) that the threshold $r = 0$

$$\mathbf{w}^T \mathbf{x} + b \geq r \iff \underbrace{\mathbf{w}^T \mathbf{x} + b - r}_{\triangleq w_0} \geq 0$$

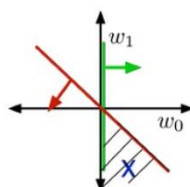
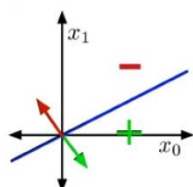
- **Eliminating the bias:** Add a dummy feature x_0 which always takes the value 1. the weight $w_0 = b$ is equivalent to a bias (same as linear regression)
- Simplified model: receive input $\mathbf{x} \in \mathbb{R}^{D+1}$ with $x_0 = 1$

$$z = \mathbf{w}^T \mathbf{x}$$

- **Example:**

x_0	x_1	t
1	0	1
1	1	0

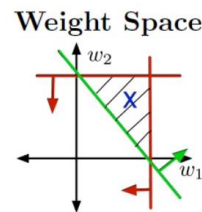
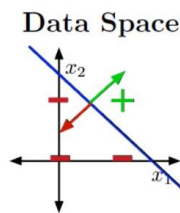
- Suppose this is our training set, with the dummy feature x_0 included.
- Which conditions on w_0, w_1 guarantee perfect classification?
 - When $x_1 = 0$, need: $z = w_0 x_0 + w_1 x_1 \geq 0 \iff w_0 \geq 0$
 - When $x_1 = 1$, need: $z = w_0 x_0 + w_1 x_1 < 0 \iff w_0 + w_1 < 0$
 - Example solution: $w_0 = 1, w_1 = -2$



$$\begin{aligned} w_0 &\geq 0 \\ w_0 + w_1 &< 0 \end{aligned}$$

- Training examples are points
- Weights (hypotheses) \mathbf{w} can be represented by **half-spaces**: $H^+ = \{\mathbf{x}: \mathbf{w}^T \mathbf{x} \geq 0\}$, $H^- = \{\mathbf{x}: \mathbf{w}^T \mathbf{x} < 0\}$
- The boundary is the **decision boundary**: $\{\mathbf{x}: \mathbf{w}^T \mathbf{x} = 0\}$
- If the training examples can be perfectly separated by a linear decision rule, we say **data is linearly separable**.
- Weights (hypotheses) \mathbf{w} are points
- Each training example \mathbf{x} specifies a half-space \mathbf{w} must lie in to be correctly classified: $\mathbf{w}^T \mathbf{x} \geq 0$ if $t = 1$
 - $x_0 = 1, x_1 = 0, t = 1 \implies (w_0, w_1) \in \mathbf{w}: w_0 \geq 0$
 - $x_0 = 1, x_1 = 1, t = 0 \implies (w_0, w_1) \in \mathbf{w}: w_0 + w_1 < 0$
- The region satisfying all the constraints is the feasible region; if this region is nonempty, the problem is

feasible, otherwise it is infeasible.



- Slice for $x_0 = 1$ and
- Example sol: $w_0 = -1.5, w_1 = 1, w_2 = 1$
- Decision boundary:

$$w_0x_0 + w_1x_1 + w_2x_2 = 0$$

$$\implies -1.5 + x_1 + x_2 = 0$$
- Slice for $w_0 = -1.5$ for the constraints
- $w_0 < 0$
- $w_0 + w_2 < 0$
- $w_0 + w_1 < 0$
- $w_0 + w_1 + w_2 \geq 0$

3.2 Towards Logistic Regression

Define loss function then try to minimize the resulting cost function

Attempt 1: 0-1 loss

$$\mathcal{L}_{0-1}(y, t) = \begin{cases} 0 & \text{if } y = t \\ 1 & \text{if } y \neq t \end{cases}$$

$$= \mathbb{I}[y \neq t]$$

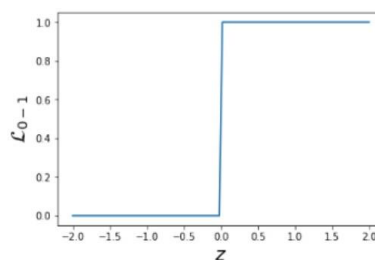
- Usually, the cost \mathcal{J} is the averaged loss over training examples; for 0-1 loss, this is the misclassification rate:

$$\mathcal{J} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[y^{(i)} \neq t^{(i)}]$$

- Minimum of a function will be at its critical points, use **Chain rule** to find the critical point of 0-1 loss

$$\frac{\partial \mathcal{L}_{0-1}}{\partial w_j} = \frac{\partial \mathcal{L}_{0-1}}{\partial z} \frac{\partial z}{\partial w_j}$$

- $\partial \mathcal{L}_{0-1} / \partial z$ is zero everywhere it's defined:



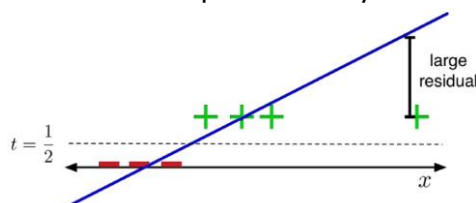
- $\partial \mathcal{L}_{0-1} / \partial w_j = 0$ means that changing the weights by a very small amount probably has no effect on the loss.
- Almost any point has 0 gradient!

Attempt 2: Linear Regression

$$z = \mathbf{w}^T \mathbf{x}$$

$$\mathcal{L}_{SE}(z, t) = \frac{1}{2} (z - t)^2$$

- Doesn't matter that the targets are actually binary. Treat them as continuous values.
- For this loss function, it makes sense to make final predictions by thresholding z at $1/2$



- The loss function hates when you make correct predictions with high confidence!

- It $t = 1$, it's more unhappy about $z = 10$ than $z = 0$.

Attempt 3: Logistic Activation Function

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

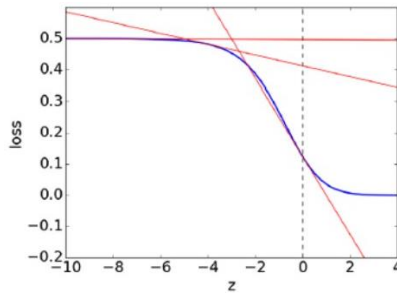
- $\sigma^{-1}(y) = \log(y/(1 - y))$ is called the **logit**.
- A linear model with a logistic nonlinearity is known as **log-linear**:

$$z = \mathbf{w}^T \mathbf{x}$$

$$y = \sigma(z)$$

$$\mathcal{L}_{SE}(y, t) = \frac{1}{2} (y - t)^2$$

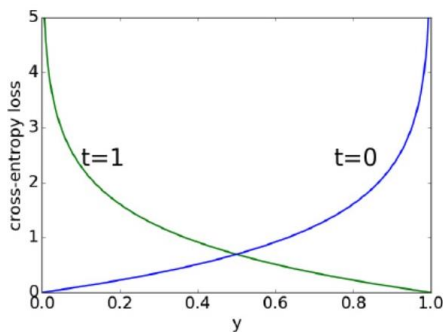
- Used in this way, σ is called an **activation function**.



(plot of \mathcal{L}_{SE} as a function of z , assuming $t = 1$)

$$\frac{\partial \mathcal{L}}{\partial w_j} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial w_j}$$

- For $z \ll 0$, we have $\sigma(z) \approx 0$.
- $\partial \mathcal{L} / \partial z \approx 0$ (check) $\rightarrow \partial \mathcal{L} / \partial w_j \approx 0 \rightarrow$ derivative w.r.t. w_j is small $\rightarrow w_j$ is like a critical point.
- If the prediction is really wrong, you should be far from a critical point (which is your candidate solution).
- Because $y \in [0, 1]$, we can interpret it as the estimated probability that $t = 1$. If $t = 0$, then we want to heavily penalize $y \approx 1$.
- Cross-entropy loss (aka log loss) captures this intuition:



$$\begin{aligned} \mathcal{L}_{CE}(y, t) &= \begin{cases} -\log y & \text{if } t = 1 \\ -\log(1 - y) & \text{if } t = 0 \end{cases} \\ &= -t \log y - (1 - t) \log(1 - y) \end{aligned}$$

- The logistic loss is a **convex function** in \mathbf{w} , so let's consider the **gradient descent** method.
 - Recall: we initialize the weights to something reasonable and repeatedly adjust them in the direction of steepest descent.
 - A standard initialization is $\mathbf{w} = 0$.

$$\mathcal{L}_{CE}(y, t) = -t \log y - (1 - t) \log(1 - y)$$

$$y = 1 / (1 + e^{-z}) \text{ and } z = \mathbf{w}^T \mathbf{x}$$

$$\begin{aligned} \frac{\partial \mathcal{L}_{CE}}{\partial w_j} &= \frac{\partial \mathcal{L}_{CE}}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w_j} = \left(-\frac{t}{y} + \frac{1-t}{1-y} \right) \cdot y(1-y) \cdot x_j \\ &= (y - t) x_j \end{aligned}$$

Gradient descent (coordinate-wise) update to find the weights of logistic regression:

$$\begin{aligned}
 w_j &\leftarrow w_j - \alpha \frac{\partial \mathcal{J}}{\partial w_j} \\
 &= w_j - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)}
 \end{aligned}$$

Gradient descent updates for Linear regression and Logistic regression (both examples of generalized linear models):

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) \mathbf{x}^{(i)}$$

3.3 Multiclass Classification and Softmax Regression

3.3.1 Multiclass Classification

- Classification tasks with more than two categories
- Targets form a discrete set $\{1, \dots, K\}$.
- It's often more convenient to represent them as **one-hot vectors**, or a **one-of-K encoding**:

$$\mathbf{t} = (0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^K$$

$\underbrace{\hspace{10em}}$
 entry k is 1

- We can start with a linear function of the inputs.

$$z_k = \sum_{j=1}^D w_{kj} x_j + b_k \quad \text{for } k = 1, 2, \dots, K$$

- Now there are D input dimensions and K output dimensions, so we need $K \times D$ weights, which we arrange as a weight matrix \mathbf{W} .
- Also, we have a K -dimensional vector \mathbf{b} of biases. Then eliminate the bias \mathbf{b} by taking $\mathbf{W} \in \mathbb{R}^{K \times (D+1)}$ and adding a dummy variable $x_0 = 1$.

$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b} \quad \text{or with dummy } x_0 = 1 \quad \mathbf{z} = \mathbf{W}\mathbf{x}$$

- We can interpret the magnitude of z_k as a measure of how much the model prefers k as its prediction to turn this linear prediction into a one-hot prediction.

$$y_i = \begin{cases} 1 & i = \arg \max_k z_k \\ 0 & \text{otherwise} \end{cases}$$

3.3.2 Softmax Regression

- We need to soften our predictions for the sake of optimization.
- We want soft predictions that are like probabilities, i.e., $0 \leq y_k \leq 1$ and $\sum_k y_k = 1$.
- A natural activation function to use is the softmax function, a multivariable generalization of the logistic function:

$$y_k = \text{softmax}(z_1, \dots, z_K)_k = \frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}}$$

- Outputs can be interpreted as probabilities (positive and sum to 1)
- If z_k is much larger than the others, then $\text{softmax}(\mathbf{z})_k \approx 1$ and it behaves like $\arg \max$.

- If a model outputs a vector of class probabilities, we can use cross-entropy as the loss function:

$$\begin{aligned}
 \mathcal{L}_{\text{CE}}(\mathbf{y}, \mathbf{t}) &= - \sum_{k=1}^K t_k \log y_k \\
 &= -\mathbf{t}^T (\log \mathbf{y})
 \end{aligned}$$

- Just like with logistic regression, we typically combine the softmax and cross-entropy into a **softmax-cross-entropy function**.
- Softmax regression (with dummy $x_0 = 1$):

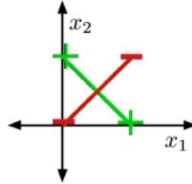
$$\begin{aligned}
 \mathbf{z} &= \mathbf{W}\mathbf{x} \\
 \mathbf{y} &= \text{softmax}(\mathbf{z}) \\
 \mathcal{L}_{\text{CE}} &= -\mathbf{t}^T (\log \mathbf{y})
 \end{aligned}$$

- Gradient descent updates can be derived for each row of \mathbf{W} :

$$\frac{\partial \mathcal{L}_{\text{CE}}}{\partial \mathbf{w}_k} = \frac{\partial \mathcal{L}_{\text{CE}}}{\partial z_k} \cdot \frac{\partial z_k}{\partial \mathbf{w}_k} = (y_k - t_k) \cdot \mathbf{x}$$

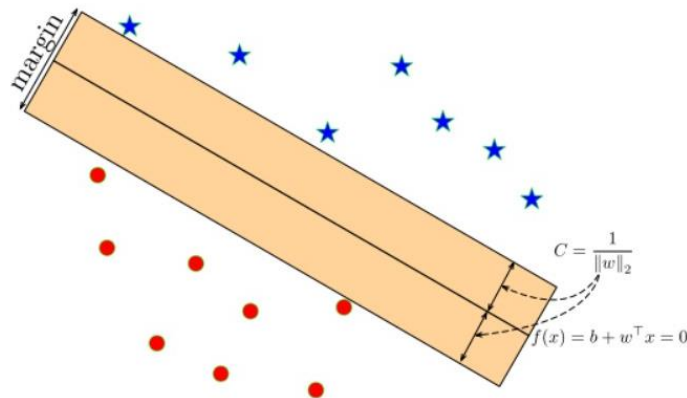
$$\mathbf{w}_k \leftarrow \mathbf{w}_k - \alpha \frac{1}{N} \sum_{i=1}^N (y_k^{(i)} - t_k^{(i)}) \mathbf{x}^{(i)}$$

- Similar to linear/logistic reg (no coincidence) (verify the update)
- Sometimes we can overcome the nonlinear problem with feature maps:

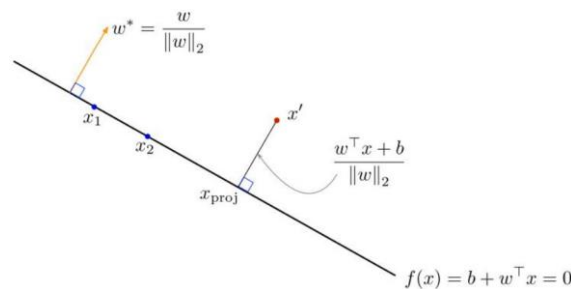


4 SVM, SVM Loss and Softmax Loss

4.1 Optimal Separating Hyperplane



- Concept: A hyperplane that separates two classes and maximizes the distance to the closest point from either class, i.e., maximize the margin of the classifier.
- Intuitively, ensuring that a classifier is not too close to any data points leads to better generalization on the test data.



4.1.1 Geometry of Points and Planes

- Recall that the decision hyperplane is orthogonal (perpendicular) to \mathbf{w} .
- The vector $\mathbf{w}^* = \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$ is a unit vector pointing in the same direction as \mathbf{w} .
- The same hyperplane could equivalently be defined in terms of \mathbf{w}^*
- The (signed) distance of a point \mathbf{x}' to the hyperplane is:

$$\frac{\mathbf{w}^T \mathbf{x}' + b}{\|\mathbf{w}\|_2}$$

4.1.2 Maximizing Margin as an Optimization Problem

- Recall: the classification for the i -th data point is correct when:
- $$\text{sign}(\mathbf{w}^T \mathbf{x}^{(i)} + b) = t^{(i)}$$

- This can be rewritten as:

$$t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) > 0$$

- Enforcing a margin of C:

$$t^{(i)} \cdot \underbrace{\frac{(\mathbf{w}^\top \mathbf{x}^{(i)} + b)}{\|\mathbf{w}\|_2}}_{\text{signed distance}} \geq C$$

- Max-margin objective:

$$\begin{aligned} \max_{\mathbf{w}, b} \quad & C \\ \text{s.t.} \quad & \frac{t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b)}{\|\mathbf{w}\|_2} \geq C \quad i = 1, \dots, N \end{aligned}$$

- Plug in $C = 1/\|\mathbf{w}\|_2$ and simplify:

$$\underbrace{\frac{t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b)}{\|\mathbf{w}\|_2} \geq \frac{1}{\|\mathbf{w}\|_2}}_{\text{geometric margin constraint}} \iff \underbrace{t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1}_{\text{algebraic margin constraint}}$$

- Equivalent optimization objective:

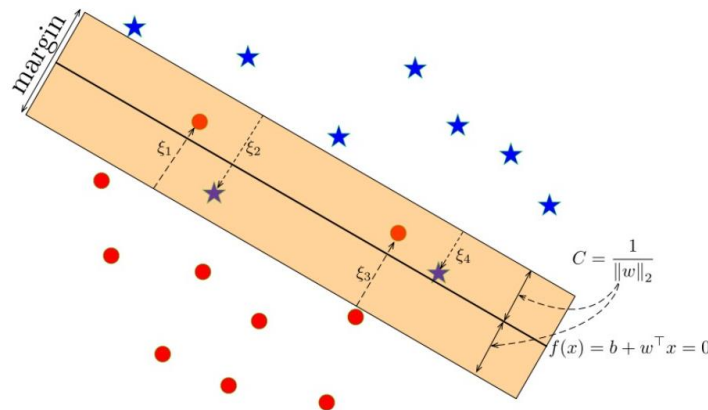
$$\begin{aligned} \min \quad & \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1 \quad i = 1, \dots, N \end{aligned}$$

4.2 Support Vector Machine

4.2.1 Concepts

- Observe: if the margin constraint is not tight for $\mathbf{x}^{(i)}$, we could remove it from the training set and the optimal \mathbf{w} would be the same.
- The important training examples are the ones with algebraic margin 1, and are called **support vectors**
- Hence, this algorithm is called the (hard) **Support Vector Machine (SVM)** (or Support Vector Classifier).
- SVM-like algorithms are often called **max-margin** or **large-margin**.

4.2.2 Maximizing Margin for Non-Separable Data Points



- Main idea:
 - Allow some points to be within the margin or even be misclassified; were present this with **slack variables** ξ_i .
 - But constrain or penalize the total amount of slack.
- Soft margin constraint:

$$\text{For } \xi_i \geq 0: \quad \frac{t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b)}{\|\mathbf{w}\|_2} \geq C(1 - \xi_i)$$

- Reduce ξ_i
- **Soft-margin SVM objective:**

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + \gamma \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1 - \xi_i \quad i = 1, \dots, N \\ & \xi_i \geq 0 \quad i = 1, \dots, N \end{aligned}$$

- γ is a hyperparameter that trades off the margin with the amount of slack
 - For $\gamma = 0$, we'll get $\mathbf{w} = 0$.
 - As $\gamma \rightarrow \infty$ we get the hard-margin objective.
- Note: it is also possible to constrain $\sum_i \xi_i$ instead of penalizing it.

4.2.3 From Margin Violation to Hinge Loss

Let's simplify the soft margin constraint by eliminating ξ_i . Recall:

$$\begin{aligned} t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) &\geq 1 - \xi_i & i = 1, \dots, N \\ \xi_i &\geq 0 & i = 1, \dots, N \end{aligned}$$

- Rewrite as $\xi_i \geq 1 - t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b)$.
- **Case 1:** $1 - t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \leq 0$
 - ▶ The smallest non-negative ξ_i that satisfies the constraint is $\xi_i = 0$.
- **Case 2:** $1 - t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) > 0$
 - ▶ The smallest ξ_i that satisfies the constraint is $\xi_i = 1 - t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b)$.
- Hence, $\xi_i = \max\{0, 1 - t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b)\}$.
- Therefore, the slack penalty can be written as

$$\sum_{i=1}^N \xi_i = \sum_{i=1}^N \max\{0, 1 - t^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b)\}.$$

If we write $y^{(i)}(\mathbf{w}, b) = \mathbf{w}^\top \mathbf{x} + b$, then the optimization problem can be written as

$$\min_{\mathbf{w}, b, \xi} \sum_{i=1}^N \max\{0, 1 - t^{(i)} y^{(i)}(\mathbf{w}, b)\} + \frac{1}{2\gamma} \|\mathbf{w}\|_2^2$$

- The loss function $\mathcal{L}_H(y, t) = \max\{0, 1 - ty\}$ is called the **hinge loss**.
- The second term is the L_2 -norm of the weights.
- Hence, the soft-margin SVM can be seen as a linear classifier with hinge loss and an L_2 regularizer.

4.2.4 Multiclass SVM Loss

Suppose: 3 training examples, 3 classes.

With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	10.9

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$\begin{aligned} L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\ &= \max(0, 2.2 - (-3.1) + 1) \\ &\quad + \max(0, 2.5 - (-3.1) + 1) \\ &= \max(0, 5.3) + \max(0, 5.6) \\ &= 5.3 + 5.6 \\ &= 10.9 \end{aligned}$$

Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

and the full training loss is the mean over all examples in the training data:

$$L = \frac{1}{N} \sum_{i=1}^N L_i$$

$$\begin{aligned} L &= (2.9 + 0 + 10.9)/3 \\ &= 4.6 \end{aligned}$$

4.3 Softmax

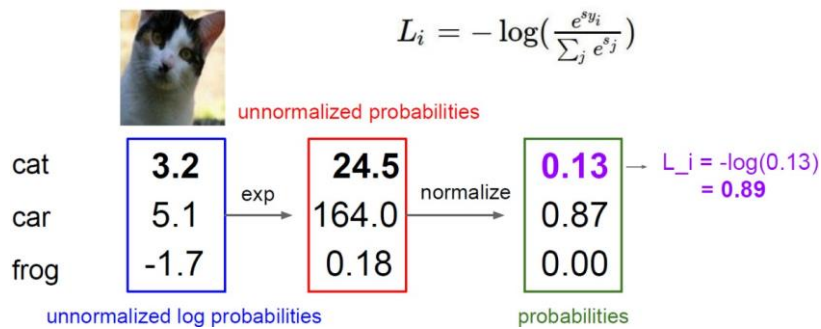
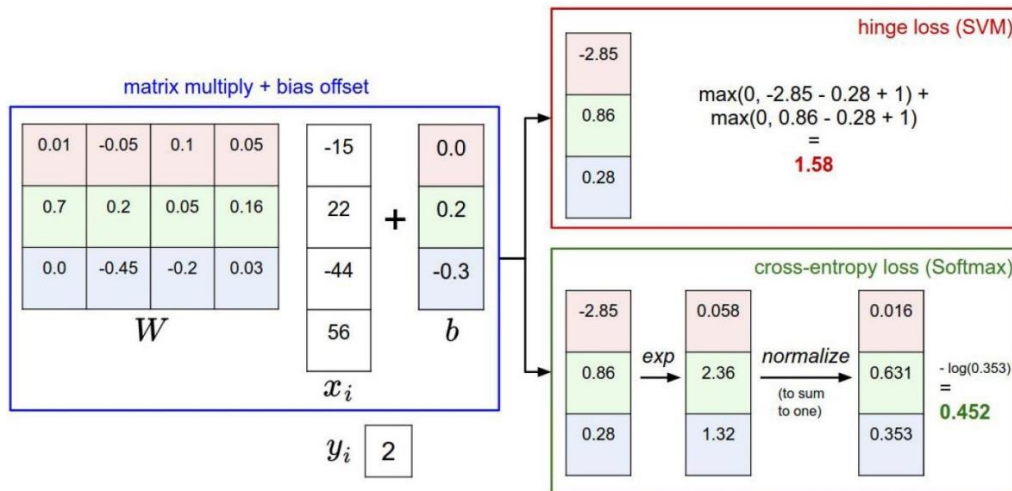
scores = unnormalized log probabilities of the classes.

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{where} \quad s = f(x_i; W)$$

Want to maximize the log likelihood, or (for a loss function) to minimize the negative log likelihood of the correct class:

$$L_i = -\log P(Y = y_i | X = x_i)$$

Softmax function

Softmax Classifier (Multinomial Logistic Regression)**4.4 SVM & Softmax**
[【深度学习 CV】SVM, Softmax 损失函数_svm 评估函数-CSDN 博客](#)
**Softmax vs. SVM**

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

assume scores:
 [10, -2, 3]
 [10, 9, 9]
 [10, -100, -100]
 and $y_i = 0$

Q: Suppose I take a datapoint and I jiggle a bit (changing its score slightly). What happens to the loss in both cases?

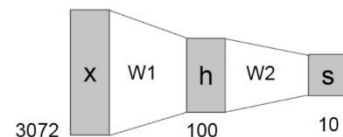
5 Neural Network and Back Propagation**5.1 Neural Network**

Neural Network: without the brain stuff

(Before) Linear score function: $f = Wx$

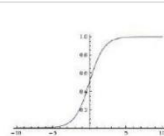
(Now) 2-layer Neural Network: $f = W_2 \max(0, W_1 x)$

or 3-layer Neural Network: $f = W_3 \max(0, W_2 \max(0, W_1 x))$

**5.1.1 Activation Functions**

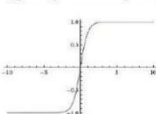
Sigmoid

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



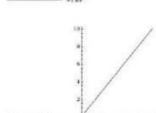
tanh

$$\tanh(x)$$



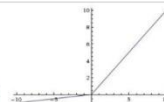
ReLU

$$\max(0, x)$$



Leaky ReLU

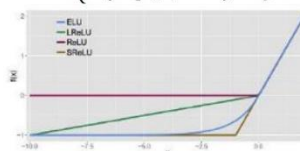
$$\max(0.1x, x)$$



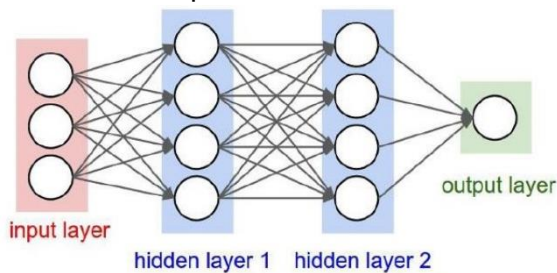
$$\text{Maxout } \max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$



Example Feed-forward computation of a Neural Network



```
# forward-pass of a 3-layer neural network:
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)
out = np.dot(W3, h2) + b3 # output neuron (1x1)
```

5.1.2 Gradient Descent

$$s = f(x; W) = Wx$$

scores function

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

SVM loss

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2$$

data loss + regularization

want $\nabla_W L$ Convolutional Network
(AlexNet)

Neural Turing Machine

input image
weights

loss

input tape

loss

Example 1:

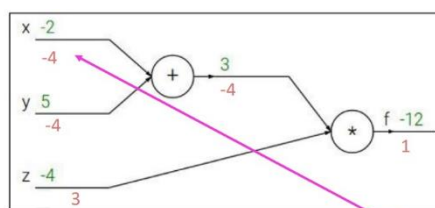
$$f(x, y, z) = (x + y)z$$

$$\text{e.g. } x = -2, y = 5, z = -4$$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

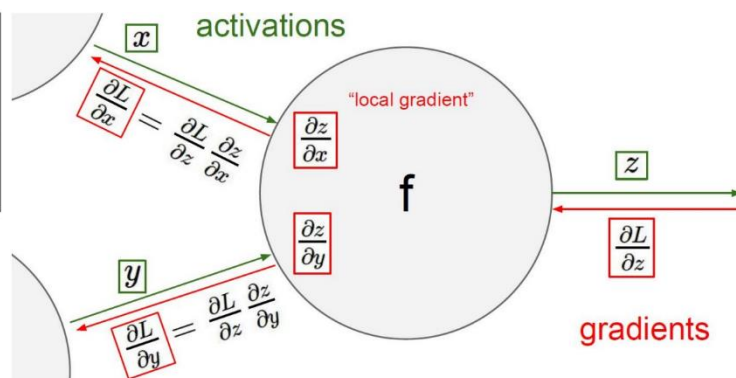
$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\text{Want: } \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$



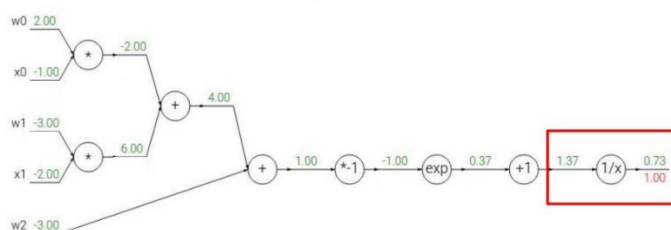
Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$



Example 2:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2 x_2)}}$$



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

→

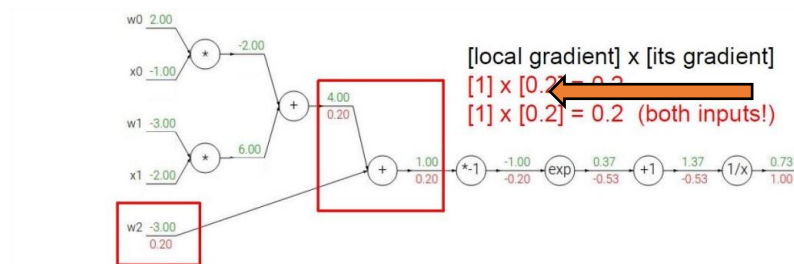
$$\frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = 1$$

求导



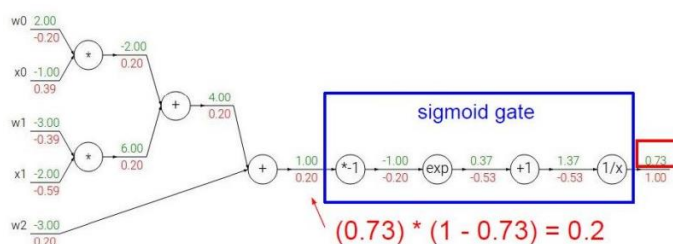
Sigmoid feature:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2 x_2)}}$$

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

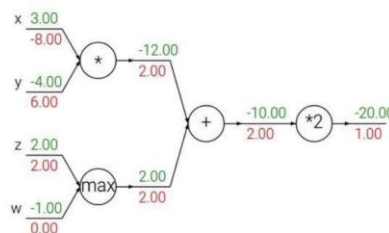
sigmoid function

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left(\frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x)) \sigma(x)$$



Patterns in backward flow

add gate: gradient distributor
 max gate: gradient router
 mul gate: gradient... "switcher"?

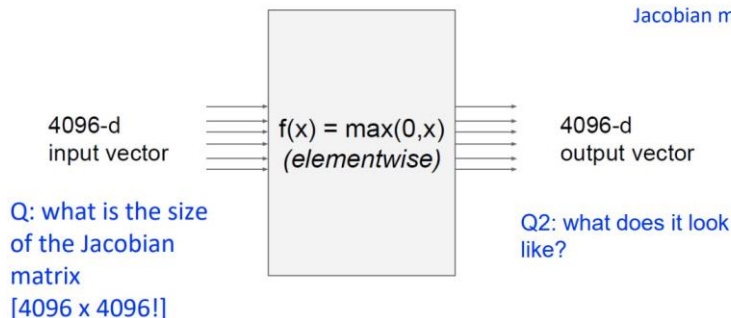


5.1.3 Gradients for vector

Vectorized operation

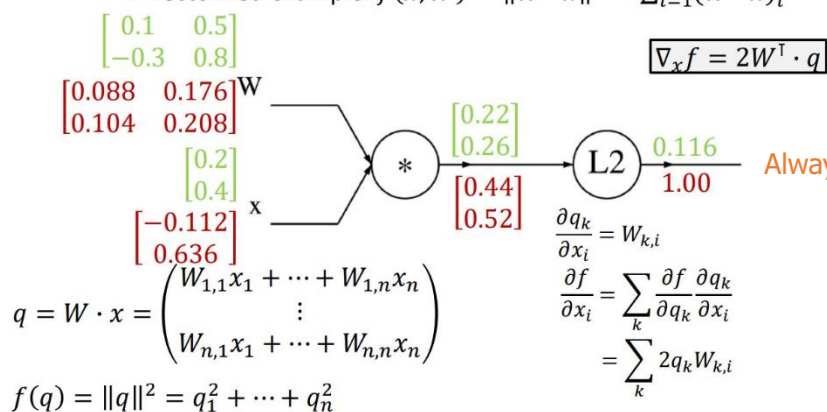
$$\frac{\partial L}{\partial x} = \left[\frac{\partial f}{\partial x} \right] \frac{\partial L}{\partial f}$$

Jacobian matrix



Example:

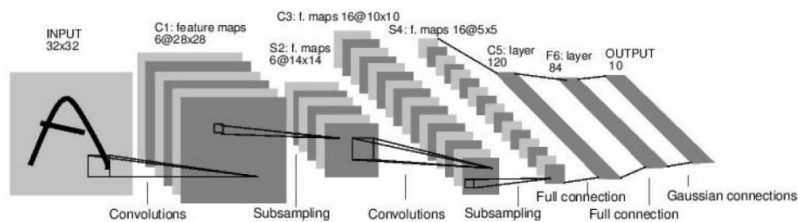
A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



Always check: The gradient with respect to a variable should have the same shape as the variable

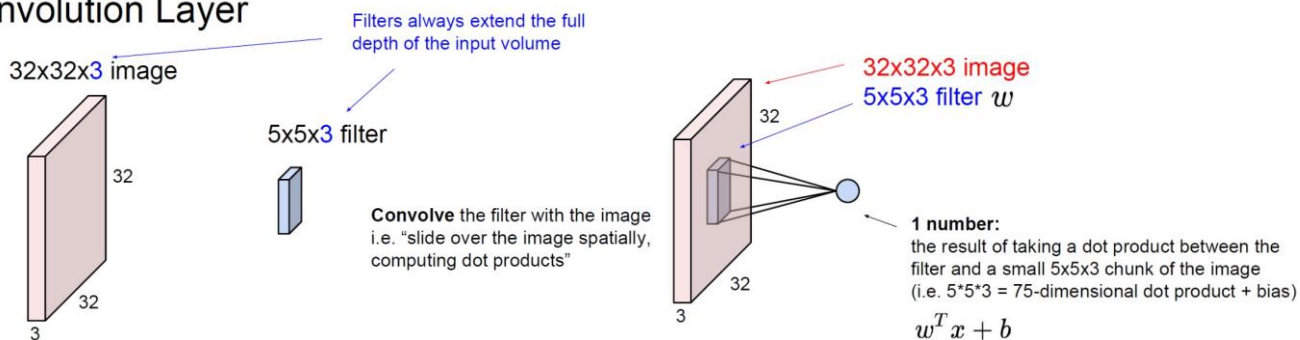
6 Convolutional Neural Network

6.1 Basic concepts

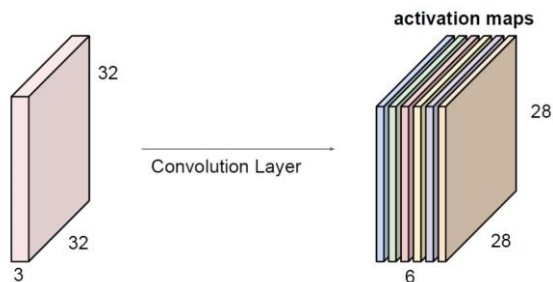


[LeNet-5, LeCun 1980]

Convolution Layer

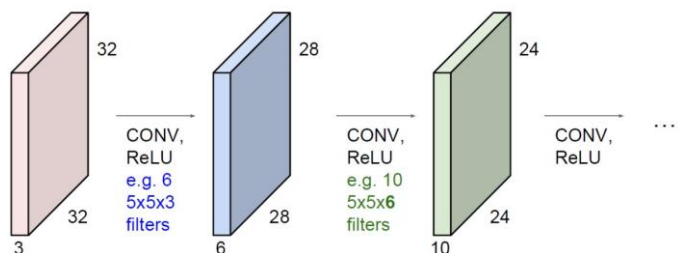


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

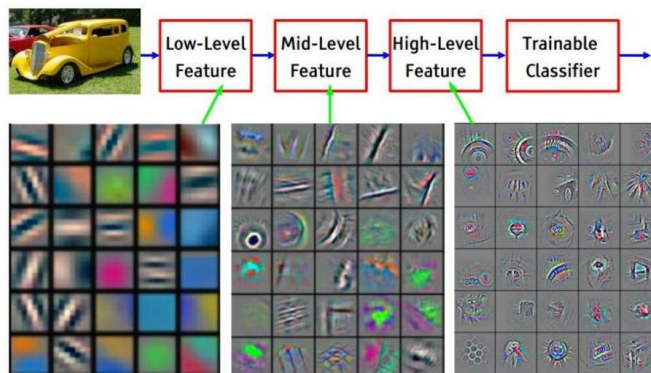


We stack these up to get a "new image" of size 28x28x6!

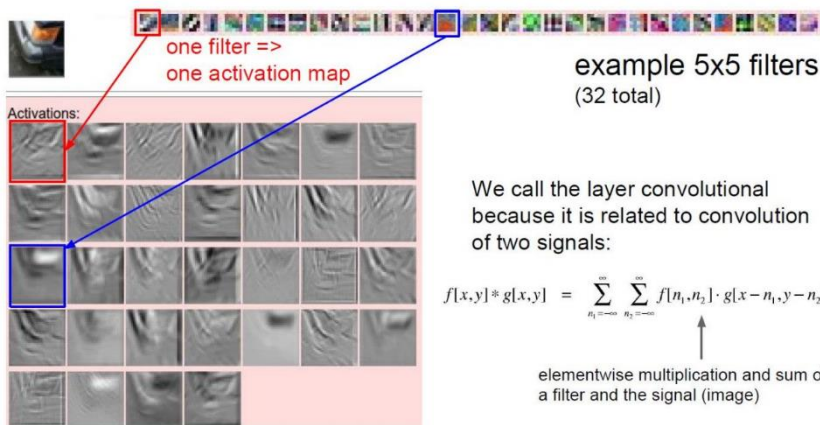
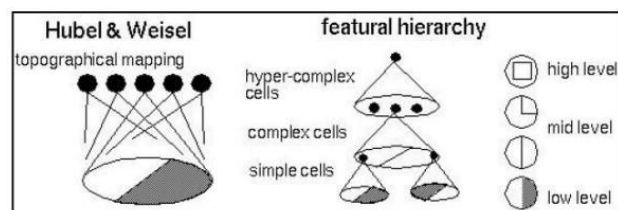
Preview: ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



Preview:

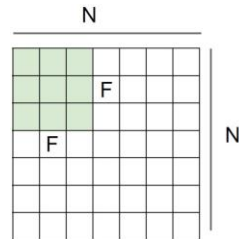


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



6.2 Convolutional layer

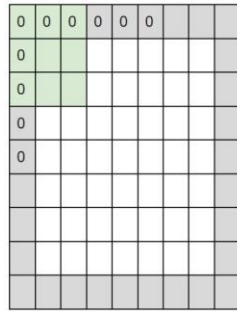
- Output size:



Output size:
 $(N - F) / \text{stride} + 1$

e.g. $N = 7, F = 3$:
 stride 1 $\Rightarrow (7 - 3) / 1 + 1 = 5$
 stride 2 $\Rightarrow (7 - 3) / 2 + 1 = 3$
 stride 3 $\Rightarrow (7 - 3) / 3 + 1 = 2.33$

- Common to zero pad the border:



e.g. input 7x7
3x3 filter, applied with **stride 1**
pad with 1 pixel border \Rightarrow what is the output?

7x7 output!

in general, common to see CONV layers with stride 1, filters of size $F \times F$, and zero-padding with $(F-1)/2$. (will preserve size spatially)

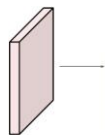
e.g. $F = 3 \Rightarrow$ zero pad with 1
 $F = 5 \Rightarrow$ zero pad with 2
 $F = 7 \Rightarrow$ zero pad with 3

- Calculation example:

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

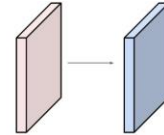
Output volume size:
 $(32 + 2 \cdot 2 - 5) / 1 + 1 = 32$ spatially, so
32x32x10



Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

Number of parameters in this layer?
 each filter has $5 \cdot 5 \cdot 3 + 1 = 76$ params (+1 for bias)
 $\Rightarrow 76 \cdot 10 = 760$



- Conclusion:

Common settings:

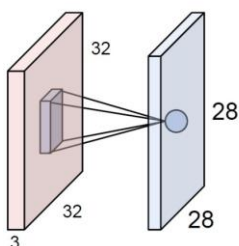
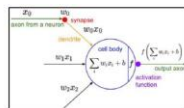
Summary. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P) / S + 1$
 - $H_2 = (H_1 - F + 2P) / S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

$K =$ (powers of 2, e.g. 32, 64, 128, 512)

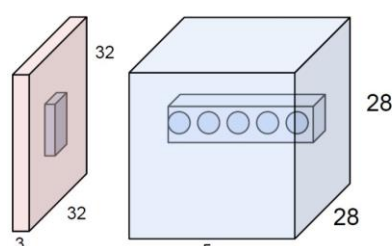
- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ?$ (whatever fits)
- $F = 1, S = 1, P = 0$

The brain/neuron view of CONV Layer



An activation map is a 28x28 sheet of neuron outputs:
 1. Each is connected to a small region in the input
 2. All of them share parameters

"5x5 filter" \rightarrow "5x5 receptive field for each neuron"

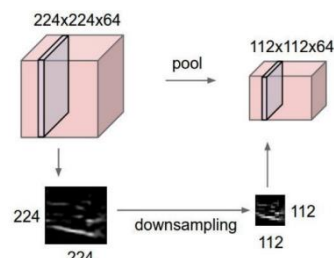


E.g. with 5 filters, CONV layer consists of neurons arranged in a 3D grid (28x28x5)

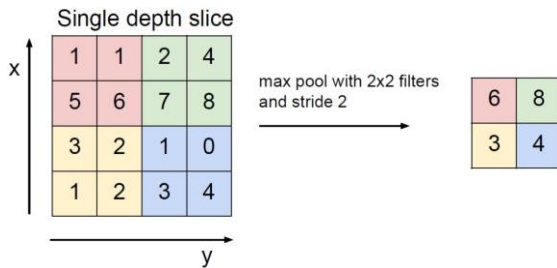
There will be 5 different neurons all looking at the same region in the input volume

6.3 Pooling layer

- Makes the representations smaller and more manageable
- Operates over each activation map independently



Max pooling:



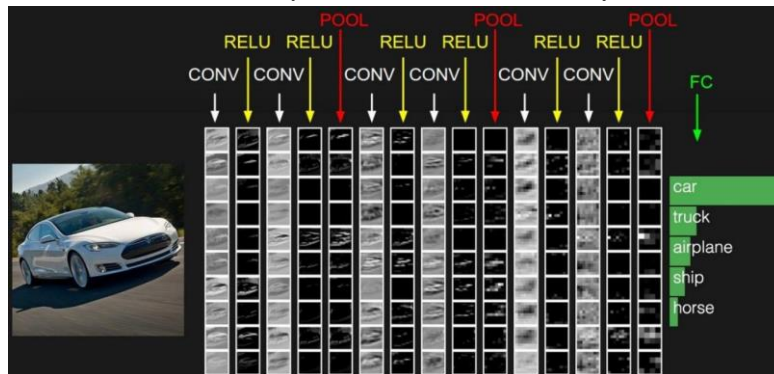
- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
 - their spatial extent F ,
 - the stride S ,
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F) / S + 1$
 - $H_2 = (H_1 - F) / S + 1$
 - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

Common settings:

 $F = 2, S = 2$ $F = 3, S = 2$

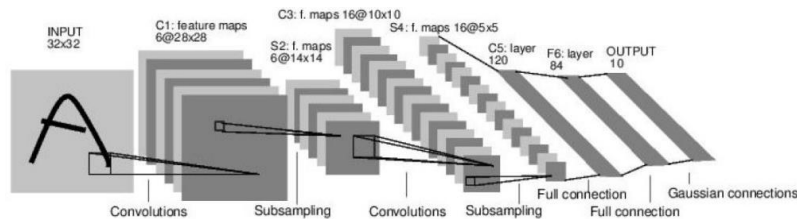
6.4 Fully connected layer

- Contains neurons that connect to the entire input volume, as in ordinary Neural Network

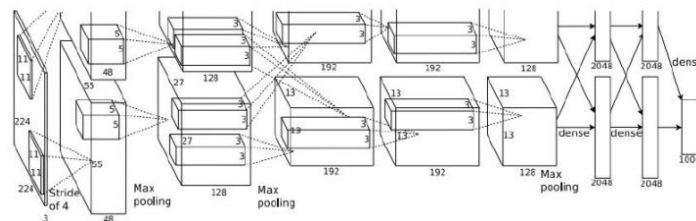


6.5 Case study: Models

6.5.1 LeNet-5



6.5.2 AlexNet



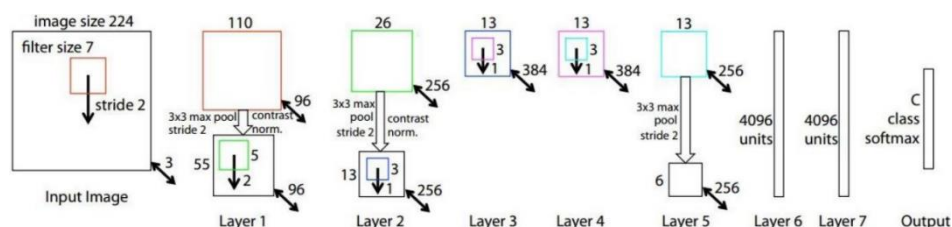
Full (simplified) AlexNet architecture:

[227x227x3] INPUT
 [55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
 [27x27x96] MAX POOL1: 3x3 filters at stride 2
 [27x27x96] NORM1: Normalization layer
 [27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
 [13x13x256] MAX POOL2: 3x3 filters at stride 2
 [13x13x256] NORM2: Normalization layer
 [13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
 [13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
 [13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
 [6x6x256] MAX POOL3: 3x3 filters at stride 2
 [4096] FC6: 4096 neurons
 [4096] FC7: 4096 neurons
 [1000] FC8: 1000 neurons (class scores)

Details/Retrospectives:

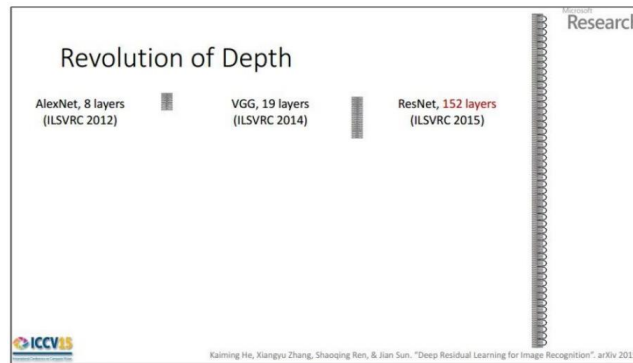
- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate $1e-2$, reduced by 10 manually when val accuracy plateaus
- L2 weight decay $5e-4$
- 7 CNN ensemble: 18.2% \rightarrow 15.4%

6.5.3 ZFNet



6.5.6 ResNet

ILSVRC 2015 winner (3.6% top 5 error)



Residual block

