



gd

Gabriele Dragotto
gabriele.dragotto@polymtl.ca

FITTING, REG AND TRAINING

OVERFITTING + VARIANCE - BIAS (SHOULD BALANCE)
UNDERFITTING + BIAS - VARIANCE

$$\mathbb{E}[\hat{y} - \tilde{f}(x)] = \mathbb{E}[(\mathbb{E}(\tilde{f}(x)) - \tilde{f}(x))^2] + \mathbb{E}(\tilde{f}(x)^2) - \mathbb{E}(\tilde{f}(x))^2 + \sigma^2$$

$$ACC = \frac{\# CORRECT CLASS}{\# SAMPLES} \quad // \text{SENS/RECALL} = \frac{\# TRUE POS}{\# POSITIVES} \quad ER = \frac{\# FP + \# FN}{\# SAMPLES}$$

$$SPEC = \frac{\# TRUE NEG}{\# NEG} \quad // \text{PRECISION} = \frac{\# TRUE POS}{\# DECL POS} \quad F_1 = \frac{2 PR}{P+R}$$

$$L1: |w| - \text{NON LINEAR - O-SOME} + k \cdot \text{SIGN}(w) \quad L2: \|w\|^2 - \text{SMOOTH} + 2k w \quad RM COND: 1. \sum_{k=0}^{\infty} \alpha^2(k) = \infty \quad 2. \sum_{k=0}^{\infty} \alpha^2(k) < \infty$$

LINEAR REGRESS.

$$\hat{w} = \bar{X} \bar{w} \quad EXACT \quad \frac{\partial \text{Err}(w)}{\partial w} = -2X^T(y - Xw) \quad INV: m \times m^2 \quad MULT: m^3$$

$$MSE = (\bar{y} - \bar{X} \bar{w})^T (y - Xw)$$

GENERATIVE DISCRIMINATIVE

$$P(y=1|\bar{x}) = \frac{P(\bar{x}|y=1)P(y=1)}{P(\bar{x})} \quad // \text{DISCRIMINATIVE LREG, KNN, NN} \quad // \text{GENERATIVE LDA, QDA, NB, HMM}$$

LDA $m(m+m)$ SAME \bar{Z} FOR K CLASSES, M FEATURES
NBAYES $(k-1) + m \times k$ \bar{Z} DIFF BUT DIAGONAL
QDA $m(mk+m)$ \bar{Z} DIFF

LOGISTIC REGRESSION

$$\text{LOG-ODDS} \quad a = \sum \bar{w}_i x_i = \text{LOG} \frac{p}{1-p}$$

$$\text{ODDS} \quad \sigma = b^a = b^{\sum w_i x_i}$$

$$\text{LIN. FN.} \quad e^{(\bar{w}^T \bar{x})} = \sigma(a) = \frac{1}{1+e^{-w^T x}} = \frac{1}{1+e^{-a}} \quad \text{APPROXIMATES THE BOUNDARY LINEARLY} \quad P(x_i=1) \text{ GIVEN FEATURE}$$

$$\text{W KELIT.} \quad L(D) = \prod \sigma(w^T x_i)^{x_i} (1 - \sigma(w^T x_i))^{1-x_i}$$

$$\text{LOG(L(D))} \quad \text{LOG(L(D))} = \sum x_i \text{LOG}(\sigma(w^T x_i)) + (1-x_i) \text{LOG}(1-\sigma) \quad H(D) = -\text{LOG L(D)}$$

$$\text{GDS} \quad \frac{\partial H(\bar{w})}{\partial \bar{w}} = -\sum_{i=1}^m \bar{x}_i (x_i - \sigma(\bar{w}^T \bar{x}_i)) \quad \text{WITH } \bar{w}_{k+1} = \bar{w}_k + \alpha \frac{\partial H(\bar{w})}{\partial \bar{w}} \text{ ONE } O(m) \text{ FULL } O(m^2)$$

$$\text{GAUSSIAN} \quad P(x_i|\bar{x}_i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \bar{w}^T \bar{x}_i)^2}{2\sigma^2}}$$

$$\text{LOG-L} \quad \text{LOG(L(D))} = \sum_{i=1}^m \text{LOG}(\sqrt{2\pi}\sigma^2) - \frac{(x_i - \bar{w}^T \bar{x}_i)^2}{2\sigma^2} \quad \frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z))$$

$$\text{LDA, QDA} \quad \text{COV} = \bar{\Sigma} = \sum_{k=1}^m \frac{(\bar{x}_k - \mu_k)(\bar{x}_k - \mu_k)^T}{N_0 + N_1 - 2} \rightarrow O(m^2) \quad \nabla \|x\|_2^2 = 2x$$

$$\text{LDA BOUNDARY: } \text{LOG}(\frac{P(\bar{x}|y=1)}{P(\bar{x}|y=0)}) = -\frac{1}{2} \mu_1^T \bar{\Sigma}^{-1} \mu_1 + \frac{1}{2} \mu_0^T \bar{\Sigma}^{-1} \mu_0 + x^T \bar{\Sigma}^{-1} (\mu_1 - \mu_0) \bar{w}^T \bar{x}$$

// FEATURES CONDITIONALLY INDEPENDENT

$$P(x_i|y) = P(x_j|y, x_i) \quad \forall i, j \Rightarrow P(x|y) \propto \prod_i P(x_i|y) \cdot P(y) \quad O(m) \text{ vs } O(2^m)$$

$$\Theta_i = P(y=1) \quad \text{ALWAYS ADDITIVE!}$$

$$\text{SMOOTHING: } \Theta_{j,i} = \frac{\# x_{i,j} + 1}{\# x_i + 2}$$

$$\Theta_{j,1} = P(x_j=1, y=1) \quad \Theta_{j,0} = P(x_j=1, y=0)$$

$$\text{LOG(L(D))}_{\text{LOG}} \frac{P(y=1|\bar{x})}{P(y=0|\bar{x})} = \text{LOG} \frac{P(y=1)}{P(y=0)} + \sum_{j=1}^m \frac{P(x_j|y=1)}{P(x_j|y=0)} \quad \text{SWAP } P(x_i|y_0) \rightarrow P(x_i|k) = P(k) \cdot P(x_i|k)$$

$$I(E) = \text{LOG } 1/P(E)$$

$$H(x|x) = \sum P(x=\bar{x}) H(x|x=\bar{x})$$

$$H(S) = \sum_i P_i I(S_i) = -\sum_i P_i \text{LOG} P_i$$

CLASS. MAX INFO GAIN
REG. MIN STD. DEV

FEATURES

$$TF_i = \frac{\text{OCC. WORD}_i}{\text{WORD IN DOC}} \quad \text{PCA ARGMIN} \left(\sum_{i=1}^m \|x - xWU^T\|^2 \right)$$

$$IDF_i = \text{LOG} \frac{\# \text{DOC}}{\# \text{DOC WITH } i} \quad W^{m \times m} \text{ m x m FIRST m EIGENVECTOR BY EIGENVALUE} \quad \text{ORTHOG. COLUMNS} \quad \text{1TH DIR WITH}$$

SVM

PERCEPTION $h(\bar{x}) = \text{SIGN}(\bar{w}^T \bar{x})$ ALWAYS CONVERGES IFF LINEARLY SEPARABLE

$$\text{ERR}(\bar{w}) = \sum \{0 \text{ IFF } x_i \bar{w}^T x_i \geq 0; 1 \text{ ELSE}\}$$

HARD-SVM
 $O(m^3)$

S-VECT $\bar{w}^T \bar{x} + b$
MARGIN $S = \frac{1}{2} \frac{1}{\|\bar{w}\|}$

MIN $\frac{1}{2} \|\bar{w}\|^2$

S.T. $y_i \bar{w}^T x_i \geq 1 \quad \forall i=1 \dots m$

$d_{\bar{w}, b}(x) = \frac{\bar{w}^T x + b}{\sqrt{\bar{w}^2 + b^2}}$

LAGRANGIAN

$L(\bar{w}, \alpha) = \frac{1}{2} \|\bar{w}\|^2 + \sum \alpha_i (1 - y_i (\bar{w}^T x_i))$

$M = 2 d_{\bar{w}, b}(x)$

$\frac{\partial L}{\partial \bar{w}} = \bar{w} - \sum \alpha_i x_i y_i$ ← WEIGHTS!

XOR = $x_1 x_2 + \bar{x}_1 \bar{x}_2$

OUTPUT

$h(\bar{w}) = \text{SIGN}(\bar{w}^T \bar{x} + b)$

SOFT-SVM

$L_{0-\infty} \rightarrow L_{0-1}$ **HINGE LOSS** $E = L_{\text{HIN}}(\bar{w}^T x_i, y_i) = \max[1 - \bar{w}^T x_i, 0]$

$\Rightarrow \min_{\bar{w}, \alpha} C \cdot \sum \epsilon_i + \frac{1}{2} \|\bar{w}\|^2$

$\alpha = 0$ GREAT! 😊

$\epsilon = 0$ ON MARGIN

S.T. $y_i \bar{w}^T x_i \geq 1 - \epsilon_i \quad \forall i=1 \dots m$

$\epsilon \in (0, 1)$ IN MARGIN

$\epsilon_i \geq 0 \quad \forall i$

$\epsilon = 1$ BOUNDARY

$\epsilon > 1$ MISS

$G_k = e^{-\frac{\|\bar{x} - \bar{z}_k\|^2}{2\sigma^2}}$

$S_k = \text{TANH}(C_1 x + C_2)$

KERNEL-T

MAP INTO HIGHER DIMENTIONAL SPACE

ENSEMBLE

BOOTSTRAPPING

K MODELS ON K SUB TRAINSET. -VAR + BIAS (DTREE, KNN)

RANDOM FOREST FOR (NODE) m RANDOM VAR BEST TEST

EXT. RANDOM.TR. FOR (TESTS) (m RANDOM AND RAN TEST) BEST TEST

BOOSTING
STACKING

INCREMENTALLY TRAIN -BIAS $LW = \alpha_i = \frac{1}{2} \log\left(\frac{1 - \epsilon_i}{\epsilon_i}\right)$
METACLASS ON TOP OF DIFFERENT APPROACHES

NEURALNET

FF-NN: PARAM = $I \otimes 0 + b_{bias}$

$h_i(\bar{x}) = \sigma(\bar{w}_i^T \bar{x} + b_i) \quad \forall i$

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

$\frac{\partial \sigma(\bar{x})}{\partial \bar{x}} = \sigma(\bar{x})(1 - \sigma(\bar{x}))$

$E_{\text{out}} = (\bar{y} - \bar{y})^T \bar{h}_{\text{LAST}} = \sigma_0^T \bar{h}_{\text{LAST}}$

$E_{\text{hidd}} = \sigma_{h,j}^T \bar{x}$

2 HIDDEN LAYERS
DO THE MAGIC!

SIGN, TANH SATURATE

RELU

GOOD

SOFTPLUS

SMOOTHER BUT COMP.

SOFTMAX

MULTICLASS

MOMENTUM $M = \beta \Delta_{i-1} \bar{w} \Rightarrow \Delta_i \bar{w} = \frac{\alpha}{\partial w} + M$

CNN

INVARIANT TO TRANSFORMATION. **STRIDE=SPACING**

CONV (m · m · p · k) FILTER (mxm)

P INPUT K OUTPUT

RNN

TEMPORALITY + SHARED WEIGHTS

END OF SEQ SENTIMENT **HIDDEN-S-REC** ELMAN

END OF STEP LANGUAGE **OUTPUTS-REC** JORDAN

JORDAN $\bar{h}_t = (\bar{w}_0 e_t + \bar{U} \bar{x}_t + \bar{b})$

$\bar{O}_t = \phi(\bar{V} \bar{h}_t + \bar{c})$ **BPTT** WITH GRADIENT CLIPPING OR TRUNCATED

ELMAN $\bar{h}_t = (\bar{w} \bar{h}_{t-1} + \bar{U} \bar{x}_t + \bar{b})$

$\bar{w} = \bar{Q} \bar{D} \bar{Q}^T \Rightarrow \bar{h}_t = \bar{Q} \bar{D}^d \bar{Q}^T \bar{h}_0$

LSTM

FIX VANISHING

HS-TRACK SELECT INFO FROM PAST X NEXT

CELL-STATE INFO FOR NEXTS

FORGET
INPUT
OUTPUT

GATE

AMOUNT TO KEEP
WHICH TO KEEP
WHAT TO BROADCAST

ENCODER CONTEXT VECTOR
DECODER + ATTENTION!

TEACHER FORCING FEED RIGHT ANSWER $y(t)$ AT $t+1$

MVAR NAME
BAYES

$P(x_j | y=k) = \frac{1}{\sqrt{2\pi\sigma_{j,k}^2}} e^{-\frac{(x_j - \mu_{j,k})^2}{2\sigma_{j,k}^2}}$

$\log(P(x_j | y=k)) = \log \frac{1}{\sqrt{2\pi\sigma_{j,k}^2}} - \frac{(x_j - \mu_{j,k})^2}{2\sigma_{j,k}^2} - \log(\sigma_{j,k})$

$\log\left(\frac{P(y=1|x)}{P(y=0|x)}\right) = \log(P(y=1)) + \sum_{j=1}^m \left(\frac{(x_j + \mu_{j,1})^2}{\sigma_{j,1}^2} - \log(\sigma_{j,1}) \right) + \log(P(y=0)) + \sum_{j=1}^m \left(\frac{(x_j + \mu_{j,0})^2}{\sigma_{j,0}^2} - \log(\sigma_{j,0}) \right)$

NB-NORMAL

$P(y=k | \bar{x}_i) = \frac{P(y=k) \cdot \prod_i P(x_i = \%)}{P(x_i = \%)}$

$P(x_i = \% | y=k) = P(x_i = \% | y=1) + P(x_i = \% | y=0)$

ODDS: $\frac{P(y=1 | x_i)}{P(y=0 | x_i)} = \frac{P(x_i = \% | y=1)}{P(x_i = \% | y=0)}$

$P(x_i = \% | y=k) = P(x_i = \% | y=k) \cdot P(x_i = \% | y=k)$



Canada Excellence Research Chair
Data Science for Real-Time Decision-Making

Applied Machine Learning

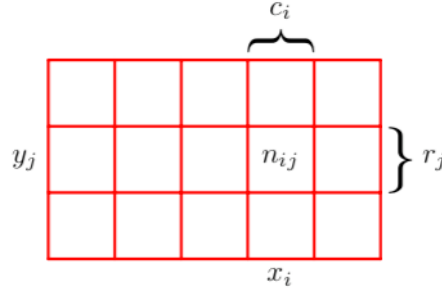
COMP551 @ McGill University

Gabriele Dragotto
gabriele.dragotto@polymtl.ca

Statistical tools

Definition 1. The *joint probability* for 2 discrete random variables X and Y is defined as $P(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$

Figure 1: Example of distribution for 2 variables. From Bishop (2006)



Definition 2. The *conditional probability* of $Y = y_j$ given $X = x_i$ is $P(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$

Definition 3. The *Bayes theorem* for 2 states that $P(A | B) = \frac{P(B|A)P(A)}{P(B)}$

Definition 4. The *sum rule* for 2 random variables X and Y states that $p(X) = \sum_Y p(X, Y)$

Definition 5. The *product rule* for 2 random variables X and Y states that $p(X, Y) = p(Y|X) \cdot P(X)$ where $p(Y|X)$ is the joint probability of Y given X .

With the Bayes Theorem, the former definition implies Equation (1)

$$p(Y|X) = \frac{p(X|Y) \cdot p(Y)}{p(X)} \quad \text{and} \quad p(X) = \sum_Y p(X|Y) \cdot p(Y) \quad (1)$$

Definition 6. The *Kolmogorov product rule* for 2 random variables X and Y states that $P(A | B) = \frac{P(A \cap B)}{P(B)}$.

If the random variable is continue, the introduced concepts cha:

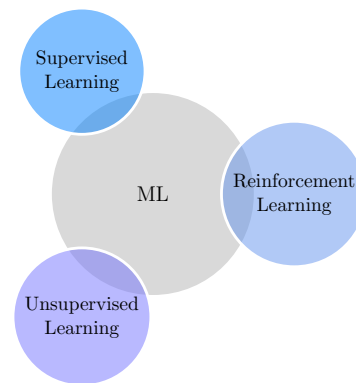
$$\int_{-\infty}^{+\infty} p(x) dx = 1 \quad p(x) \geq 0 \quad \forall x \quad (2)$$

$$P(x \leq z) = \int_{-\infty}^z p(x) dx \quad P'(x) = p(x) \quad (3)$$

1 Introduction

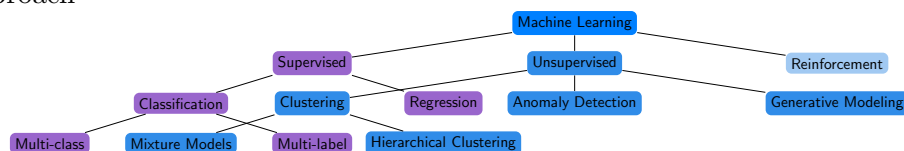
A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , **improves** with experience E

Class Notes by Hamilton (2019)



In a broader scope, *AI* can be seen as the field of Computer Science aiming to automate actions, decision and learning process typical of human beings. A sub field of *AI* is **machine learning**, which implements in machines the capability of **extracting knowledge from data** (Goodfellow et al., 2016). This approach relies heavily on the **representation** of the data it receives (eg, regression only analyses some features). The subfield of **representation learning** concerns with an approach capable of guessing the right representation to use in order to produce the output.

Finally, **deep learning** aims to build lower level representation of data, and hence builds up complex representations from simpler ones. This former approach



1.1 Fitting, regularization and validation

Definition 7. A *training set* is a formatted set of data. Each column is named *variable, feature or attribute* and each row constitutes a *training example or instance*. The desired outcome is named *target* or *output variable*.

Definition 8. An *overfitted* model is a statistical model that contains more parameters than can be justified by the data. The model has a **lower true error** on a different hypothesis than the chosen one, hence *high variance/low bias*.

Counter with: model comparison, cross-validation, regularization, early stopping, pruning, Bayesian priors, or dropout. **Reduce variance** at the cost of **some bias**.

Definition 9. *Underfitting* occurs when a statistical model or machine learning algorithm cannot adequately capture the underlying structure of the data. The model has then *low variance/high bias*.

1.2 Errors

Assuming a machine learning algorithm predicts $y = f(x) + \epsilon$ where $\epsilon \sim (\mu = 0, \sigma^2)$, then the error can be decomposed in *bias* and *variance*:

$$\text{Error} \quad \mathbb{E}[y - \tilde{f}(\bar{x})] = (\text{Bias}[\tilde{f}(\bar{x})])^2 + \text{Var}[\tilde{f}(\bar{x})] + \sigma^2 \quad (4)$$

$$\text{Bias} \quad \mathbb{E}[\tilde{f}(\bar{x})] - \tilde{f}(\bar{x}) \quad (5)$$

$$\text{Variance} \quad \mathbb{E}[\tilde{f}(\bar{x})^2] - \mathbb{E}[\tilde{f}(\bar{x})]^2 \quad (6)$$

$$(7)$$

Definition 10. *Bias* occurs when a statistical model or machine learning algorithm cannot capture relations between *input and output variables*. *underfitting*.

Definition 11. *Variance* models the sensitivity between *small changes in the input*. *overfitting*.

$$\text{Accuracy} \quad A = \frac{\# \text{Correctly Classified}}{\# \text{Samples}} \quad (8)$$

$$\text{Sensitivity/Recall} \quad R = \frac{\# \text{True Positives}}{\# \text{Positives}} \quad (9)$$

$$\text{Specificity} \quad S_c = \frac{\# \text{True Negatives}}{\# \text{Negatives}} \quad (10)$$

$$\text{Precision} \quad P = \frac{\# \text{True Positives}}{\# \text{Declared Positives}} \quad (11)$$

$$\text{Error rate} \quad E_r = \frac{\# \text{False Negatives} + \# \text{False Positives}}{\# \text{Samples}} \quad (12)$$

$$F_1 \text{ score} \quad F_1 = 2 \frac{P \cdot R}{P + R} \quad (13)$$

TruePositives	FalseNegatives
FalsePositives	TrueNegatives

Table 1: The confusion matrix structure

1.2.1 Cross Validation

Definition 12. ***k-fold** cross validation splits the training-set in k partitions, training on $k - 1$ and verifying on the remaining one. It increases the computational time by a factor of k .*

Definition 13. ***Leave-out-one** eliminates one row on the training set, and evaluates the error on it. The error estimation averages on all the iterations.*

1.3 Regularization

Definition 14. ***Lasso penalization** - or $L1$ regularization - add to the loss-function a penalty term proportional to the absolute value of estimated weights. **non linear** . Sets some coefficients to 0.*

Definition 15. ***Ridge regression** - or $L2$ regularization - add to the loss-function a penalty term proportional to the square of estimated weights. Smooths some coefficients.*

2 Linear Regression

The *I.I.D* assumption states the following:

Definition 16. A set of random variables is independent and identically distributed if each random variable has the **same probability distribution** as the others and all are **mutually independent**.

The following equations are the fundamental:

$$\text{Function} \quad f_{\bar{w}_i}(\bar{X}) = \bar{X}\bar{w}_i \quad (14)$$

$$\text{Fn. MSE Error} \quad f_{\bar{w}}^*(\bar{X}) = \arg \min_i [(\bar{y} - \bar{X}\bar{w}_i)^T(\bar{y} - \bar{X}\bar{w}_i)] \quad (15)$$

2.1 Exact approach

We minimize the gradient of the MSE representation.

$$\text{Gradient MSE} \quad \frac{\delta \text{Err}(\bar{w})}{\delta \bar{w}} = -2X^T(y - X\bar{w}) \quad (16)$$

$$\Rightarrow X^T(y - X\bar{w}) = 0 \quad (17)$$

$$\Rightarrow X^TY = X^TX\bar{w} \quad (18)$$

$$\Rightarrow \bar{\mathbf{w}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} \quad (19)$$

Operation	Cost	
Matrix Inversion	1	nm^2
Matrix multiplications	3	m^3
Total	4	$\Theta(nm^2 + m^3)$

Some **issues**:

- Numerical stability
- **Singular X**: features are dependent (no full column rank).
fix? combine, apply functions, interaction terms, etc...

2.2 Gradient descent

Weights are updated **step by step** depending on the **MSE gradient** representation. Assuming $\text{Err}(\bar{w}_0) > \text{Err}(\bar{w}_1) > \dots \text{Err}(\bar{w}_{step})$, until $|\bar{w}_{k+1} -$

$\bar{w}_k| > \epsilon$, then:

$$\text{Updated Weights} \quad \bar{w}_{k+1} = \bar{w}_k - \alpha_k \frac{\delta \text{Err}(\bar{w}_k)}{\delta \bar{w}_k} \quad (20)$$

$$\Rightarrow \bar{w}_{k+1} = \bar{w}_k + 2\alpha_k X^T(y - X\bar{w}) \quad (21)$$

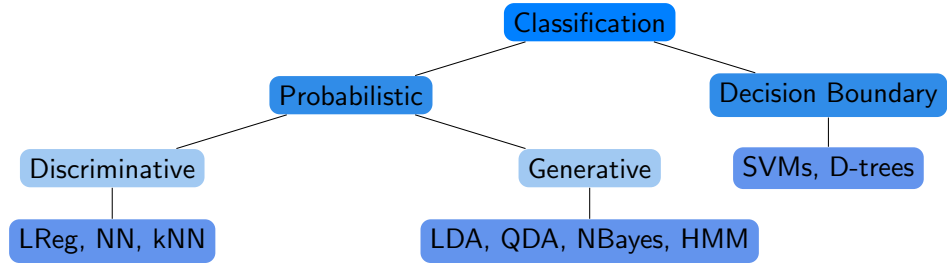
Where the function (parameter) α_k is the **learning rate** at the k -th step.

- **Large α** : gradient might **not converge** $\Rightarrow \alpha_k \rightarrow 0$ if $k \rightarrow \infty$.
- **Small α** : might loop in **local minima** $\Rightarrow \alpha_k \rightarrow 0$ if $k \rightarrow \infty$.

Claim 1. *Robbins-Monroe conditions are sufficient to ensure convergence of the \bar{w}_k to a local minimum of the error function.*

RMC: $\sum_{k=0}^{\infty} \alpha_k = \infty$ and $\sum_{k=0}^{\infty} \alpha_k^2 < \infty$

For instance, $\alpha(k) = 1/(k+1)$ satisfies the *RMC*.



$$P(y = 1|\bar{x}) = \frac{P(\bar{x}|y = 1)P(y = 1)}{P(\bar{x})}$$

Estimating **how likely are the features** given the output, and the independent **probability of the output class**.

- **Discriminative learning:** $P(y|\bar{x})$ models the **boundary** between classes.
- **Generative learning:** $P(\bar{x}|y)$ models the **distribution** of each class.

Model	Cost	Assumption
LDA	$m(m+n)$	Same covariance for k classes, and m features
Naive Bayes	$(k-1) + mk$	Different covariances but diagonal matrix
QDA	$m(mk+n)$	Different covariances for k classes

3 Logistic Regression

Definition 17. The **logit** - or log-odd ratio - is defined as the logarithm of the odds $1/(1 - p)$. Therefore $\text{logit}(p) = \log \frac{p}{1-p}$

The following equations are the fundamental:

$$\text{Log-odds} \quad a = \sum_i w_i x_i = \log \frac{p}{1-p} \quad (22)$$

$$\text{Odds} \quad o = b^a = b^{\sum_i w_i x_i} \quad (23)$$

$$\text{Lin. Logistic Fn.} \quad \sigma(\bar{w}^T \bar{x}) = \sigma(a) = \frac{1}{1 + e^{-\bar{w}^T \bar{x}}} = \frac{1}{1 + e^{-a}} \quad (24)$$

Where the **decision boundary** approximates **the log-odds with a linear function of the features** $\bar{w}^T \bar{x}$.

$\sigma(\bar{w}^T \bar{x})$ is the **probability of $y_i = 1$** given the \bar{x} input vector. The goal is to maximize the likelihood, without incurring in **numerical instability**.

$$\text{Likelihood} \quad L(D) = \prod_{i=1}^n \sigma(\bar{w}^T \bar{x})^{y_i} \cdot (1 - \sigma(\bar{w}^T \bar{x}))^{1-y_i} \quad (25)$$

$$\text{log-Likelihood} \quad \log L(D) = \sum_{i=1}^n y_i \log(\sigma(\bar{w}^T \bar{x})) + (1-y_i) \log(1-\sigma(\bar{w}^T \bar{x})) \quad (26)$$

$$\text{Cross Entropy Loss} \quad H(d) = -\log L(D) \quad (27)$$

Therefore, **maximising** likelihood corresponds to **minimizing** the cross entropy. It is easy to note that the inner term of $L(D)$:

1. $= \sigma(\bar{w}^T \bar{x})$ iff $y_i = 1$
2. $= 1 - \sigma(\bar{w}^T \bar{x})$ iff $y_i = 0$

3.1 Gradient descent

We minimize the gradient of the cross-entropy loss representation. One step is $O(nm)$, vs $O(m^3 + nm^2)$ for exact solution.

$$\text{Gradient CE} \quad \frac{\delta H(\bar{w})}{\delta \bar{w}} \quad (28)$$

$$\text{with} \quad \frac{\delta \log(\sigma)}{\delta \bar{w}} = \frac{1}{\sigma}, \quad \frac{\delta \sigma}{\delta \bar{w}} = \sigma(1 - \sigma), \quad \frac{\delta \bar{w}^T x}{\delta \bar{w}} = x, \quad (29)$$

$$\frac{\delta(1 - \sigma)}{\delta \bar{w}} = -\sigma(1 - \sigma) \quad (30)$$

$$\Rightarrow \frac{\delta \mathbf{H}(\bar{\mathbf{w}})}{\delta \bar{\mathbf{w}}} = - \sum_{i=1}^n \bar{\mathbf{x}}_i (\mathbf{y}_i - \sigma(\bar{\mathbf{w}}^T \bar{\mathbf{x}}_i)) \quad (31)$$

Therefore, the update rule

$$\text{Updated Weights} \quad \bar{w}_{k+1} = \bar{w}_k + \alpha_k \left(-\frac{\delta H(\bar{w})}{\delta \bar{w}} \right) \quad (32)$$

3.2 Probabilistic interpretation

The logistic regression can be interpreted through a gaussian distribution of the likelihood.

$$\text{Gaussian Likelihood} \quad P(y_i | \bar{x}_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \bar{w}^T \bar{x}_i)^2}{2\sigma^2}} \quad (33)$$

$$\text{Gaussian log-likelihood} \quad \log L(D) = \sum_{i=1}^n -\log(\sqrt{2\pi\sigma^2}) - \frac{(y_i - \bar{w}^T \bar{x}_i)^2}{2\sigma^2} \quad (34)$$

$$\text{Squared loss} \quad (y_i - \bar{w}^T \bar{x}_i)^2 \quad (35)$$

Therefore, maximising log-likelihood corresponds to minimizing the squared loss. Train models using theoretically grounded loss functions but evaluate using interpretable measures.

4 LDA and QDA

LDA approaches the problem by assuming that the conditional probability density functions $p(\bar{x} | y_i = \%)$ are **gaussian** functions with mean μ_i and same covariance Σ , while QDA assumes different covariances. The Bayes

optimal predicts points as being from the second class if the **log-odd** ratios is **greater than 0**.

$$P(\bar{x}|y) = \frac{1}{\sqrt{2\pi}\Sigma} e^{-\frac{(\bar{x}-\mu)^T(\bar{x}-\mu)}{2\Sigma}} \quad (36)$$

$$\text{Covariance} \quad \Sigma = \sum_{k=0}^1 \sum_{j=1}^n \frac{(\bar{x}_i - \mu_k)^T (\bar{x}_i - \mu_k)}{N_0 + N_1 - 2} \quad (37)$$

QDA has **more parameters** to estimate, but greater flexibility to estimate the target function. Estimating Σ is expensive: $O(m^2)$! The *LDA* has a linear boundary, namely $w_0 + \bar{x}^T \bar{w}$.

$$\log \frac{P(\bar{x}|y=1)}{P(\bar{x}|y=0)} = \quad (38)$$

$$\log \frac{P(y=1)}{P(x=1)} - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 + \bar{x}^T \Sigma^{-1} (\mu_1 - \mu_0) \quad (39)$$

4.1 Naive Bayes

Definition 18. The strong (naive) assumption states assumes that the features are **conditionally** independent given the output y . Therefore:

$$P(x_j|y) = P(x_j|y, x_k) \quad \forall j, k \quad \Rightarrow P(x|y) = \prod_i (x_i|y)$$

A Naive Bayes classifier with **binary** features has to estimate just $O(m)$ parameters compared to $O(2^m)$. Useful when the **number of features is high**: linear time parameters and exact form optimization. Hence, no **correlation between features** is taken into account.

$$\Theta_1 \quad P(y=1) \quad (40)$$

$$\Theta_{j,1} \quad P(x_j=1|y=1) \quad (41)$$

$$\Theta_{j,0} \quad P(x_j=1|y=0) \quad (42)$$

$$L(\Theta_1|y) \quad \Theta_1^y (1 - \Theta_1)^{1-y} \quad (43)$$

$$\text{log-Likelihood} \quad \log_L(D) = \log \frac{P(y=1|\bar{x})}{P(y=0|\bar{x})} = \quad (44)$$

$$\Rightarrow \log \frac{P(y=1)}{P(y=0)} + \log \frac{\prod_{j=1}^m P(x_j|y=1)}{\prod_{j=1}^m P(x_j|y=0)} \quad (45)$$

$$\Rightarrow \log \frac{P(y=1)}{P(y=0)} + \sum_{j=1}^m \log \frac{P(x_j|y=1)}{P(x_j|y=0)} \quad (46)$$

Therefore, **each feature contributes independently** to the classification. In order to **reduce variance**, one can exploit additive smoothing.

Definition 19. The **additive smoothing** - or Laplace smoothing - modifies estimators as: $\hat{\theta}_i = \frac{x_i + \alpha}{N + \alpha d} \quad \forall i$ where $d = p_i^{-1}$, namely the probability of class y_i .

For instance, the add one smoothing would be: $\Theta_{j,i} = \frac{\#y_i=1 \wedge x_j=1+1}{\#y_i=1+2}$.

Gaussian Naive Bayes assumes Σ is distinct between classes and diagonal, and $P(x|y)$ is assumed to be a multivariate Gaussian.

5 Decision Trees

Divide the space of features in a detailed way. Learning a tree means **learning tests** (with categorical/binary outcomes) over each branch. Usually: **grow** and **prune**.

- **Easy** to represent: Boolean functions
- **Hard** to represent: Parity function $O(2^m)$, Majority function

$$\text{Information} \quad I(E) = \log_2 \frac{1}{P(E)} \quad (47)$$

$$\text{Entropy} \quad H(S) = \sum_i p_i I(s_i) = - \sum_i p_i \log p_i \quad (48)$$

$$\text{Conditional Entropy} \quad H(y|x) = \sum_i p(x=i) H(y|x=i) \quad (49)$$

Claim 2. The entropy can be interpreted as: average amount of information per symbol, uncertainty before the outcome, number of bits for the symbol.

- **Classification**: maximize information gain.
- **Regression**: minimize standard deviation.

In order to avoid overfitting:

- **post-pruning**: prune the tree once it is fully grown.
For each node: prune iff **improves validation accuracy**. Replace decision with majority rule.
- **early-stopping**: stop the training when no information gain.

6 Feature Design

6.1 NL Features

Some language processing standard measures:

$$\text{tf} \quad tf_i = \frac{word_i}{\#words} \quad \text{in document } d_k \quad (50)$$

$$\text{idf} \quad idf_i = \log \frac{\#documents}{\#documents \text{ with } i} \quad \text{in corpus} \quad (51)$$

$$\text{tf-idf} \quad tf \cdot idf_i = tf_i \cdot idf_i \quad (52)$$

6.2 PCA

Principal component analysis - or truncated SVD - projects into a lower dimensional space a given set of features. Assume the original dimension as \mathbb{R}^m and the final one \mathbb{R}^n with $n < m$.

$$\text{PCA-problem} \quad \arg \min_{W,U} \left(\sum_i^n \|X - XWU^T\|^2 \right) \quad (53)$$

Where:

- $W^{m \times n}$: compression matrix with the first n eigenvectors sorted by **eigenvalues**. Columns are **orthogonal**, and i -th column is the i -th direction with i -th **maximal variance**.
- $U^{n \times m}$: decompression matrix.

7 Instance learning

Compute a **domain specific distance** metric between points in the training set and the ones being tested. The decision boundaries are given by the **Voronoi diagram** ran on the training data. **lazy learning**, wait for the query to generalize.

With the k -NN the prediction is the majority/mean of k -nearest points. Try to use *gaussian* distances.

8 Support Vector Machines

8.1 Perceptron

Basic linear classifier. Its error estimates how much $w^t x$ is far away from being correct. The training data is **linearly separable** iff there is **no training error**.

$$\text{Perceptron} \quad h_w(\bar{x}) = \text{sign}(\bar{w}^T \bar{x}) \quad (54)$$

$$\text{Perceptron Err} \quad \text{Err}(\bar{w}) = \sum_i \{0 \iff y_i \bar{w}^T x_i \geq 0; \text{ else } 1\} \quad (55)$$

$$(56)$$

Theorem 8.1. *If the training data is linearly separable, the perceptron will converge to 0 error in a finite number of steps. ♡ (Bishop, 2006)*

8.2 SVM

Definition 20. *A linear SVM is a perceptron with a vector \bar{w} so that the margin is maximized.*

The SVM builds the **separating hyperplane** between two classes of the input space. $\vec{w} \cdot \vec{x} - b = 0$ is the separation hyperplane.

$$\text{Separation} \quad S = \frac{2}{\|\bar{w}\|} \quad (57)$$

The optimization model, that can be trained at most in $O(n^3)$.

$$\text{minimize}_{\bar{w}} \quad \frac{1}{2} \|\bar{w}\|^2 \quad (58)$$

$$\text{subject to} \quad y_i \bar{w}^T \bar{x}_i \geq 1 \quad \forall i = 1, \dots, n \quad (59)$$

$$\|\bar{w}\| = \frac{1}{M} \quad (60)$$

The lagrangian multipliers $\bar{\alpha}$ of the dual are the **support vectors**. The weight vector is a **linear combination** of the support vectors.

Definition 21. *A support vector is a point laying on the decision boundary.*

$$\text{Lagrangian Form.} \quad L(\bar{w}, \bar{\alpha}) = \frac{1}{2} \|\bar{w}\|^2 + \sum_i \alpha_i (1 - y_i(\bar{w}^T \bar{x})) \quad (61)$$

$$\text{Derivate} \quad \frac{\delta L(\bar{w}, \bar{\alpha})}{\delta \bar{w}} = \bar{w} - \sum_i \alpha_i \bar{x}_i y_i \quad (62)$$

$$\text{Weights} \quad \bar{w} = \sum_i \alpha_i \bar{x}_i y_i \quad (63)$$

$$\text{Margin} \quad M = M(\bar{\alpha}_i) = \frac{\bar{w}^T \bar{\alpha}_i}{\|\bar{w}\|} = \frac{1}{\|\bar{w}\|} \quad (64)$$

$$\text{Output} \quad h_{\bar{w}}(\bar{x}) = \text{sign}\left(\sum_i^n \alpha_i y_i (x_i \cdot x)\right) = \text{sign}(\bar{w}^T \bar{x} + b) \quad (65)$$

The output of the classifier is given by the dotproduct of the sample x with the support vectors x_i .

8.3 Non-linearly separable class

If the training set is **not linearly separable**, either:

- **Softening** the constraint: $L_{0-\infty} \rightarrow L_{0-1}$.
- **Kernel-trick**: use non-linear kernels.

8.3.1 Soft SVM with Hinge loss

Approximate the missclassification penalty with a **linear function**.

$$\text{Hinge Loss} \quad \xi = L_{hin}(\bar{w}^T \bar{x}_i, y_i) = \max\{1 - \bar{w}^T \bar{x}_i, 0\} \quad (66)$$

Therefore optimize the following:

$$\text{minimize}_{\bar{w}, \xi} \quad C \sum_i \xi_i + \frac{1}{2} \|\bar{w}\|^2 \quad (67)$$

$$\text{subject to} \quad y_i \bar{w}^T \bar{x}_i \geq 1 - \xi_i \quad \forall i = 1, \dots, n \quad (\mathbf{Lagr} : \bar{\alpha}) \quad (68)$$

$$\xi_i \geq 0 \quad \forall i = 1, \dots, n \quad (\mathbf{Lagr} : \bar{\beta}) \quad (69)$$

With $C \rightarrow \infty$ it gives an **hard-SVM**.

- $\alpha = 0$: points outside margin.

- $\xi = 0$: points on/outside the margin.
- $\xi \in (0, 1)$: points within the margin (on with $\alpha_i > 0$).
- $\xi = 1$: points on decision line.
- $\xi > 1$: misclassified.

8.3.2 Kernel machines

The method allows to compute dot-products **without explicitly computing coordinates** in the new feature space. Optimizing with the dual requires just the computation of the [kernel function](#).

$$\text{Gaussian Kernel} \quad K = (\bar{x}, \bar{z}) = e^{-\frac{\|\bar{x} - \bar{z}\|^2}{2\sigma^2}} \quad (70)$$

$$\text{Sigmoid Kernel} \quad K = (\bar{x}, \bar{z}) = \tanh(c_1 \bar{x} \cdot \bar{z} + c_2) \quad (71)$$

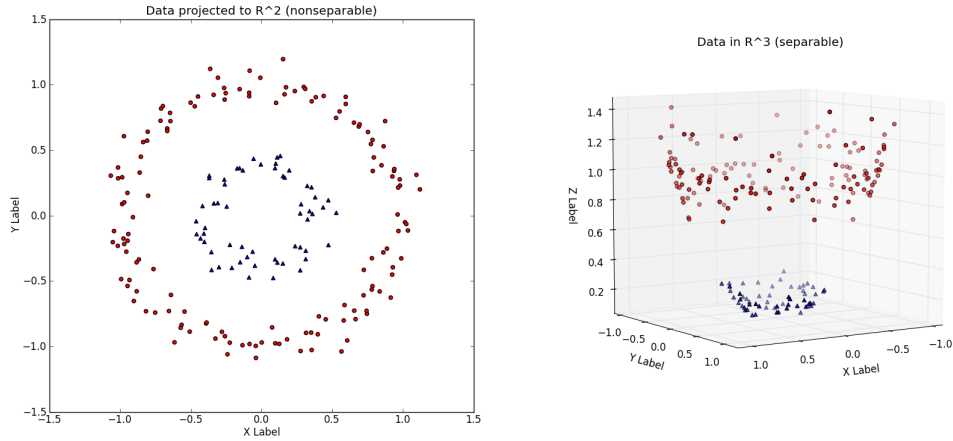


Figure 2: Kernel mapping $\mathbb{R}^2 \rightarrow \mathbb{R}^3$

9 Ensemble methods

- **Bootstrapping-Bagging**: K different models trained on *different sub training sets* (eg, sampling with replacement).
Reduce *variance* and increase *bias*: **DTree**, **kNN**.

- **Boosting**: K different models incrementally trained.
Reduce bias: **AdaBoost**. Empirically, low chance of overfitting.
misclassified are dangerous: can be weighted too much.
- **Stacking**: different classifiers combined with a meta-classifier.
Features: classifiers output.

9.1 Bagging

1 Random Forest

```

Train  $K$  different trees with  $k$  bootstraps.
for each node do
    Pick  $m$  random variables
    Determine test =  $\max\{InfoGain\}$ 
end for
Predict using ensemble

```

Each tree has high variance, but the ensemble uses averaging, which reduces variance.

2 Extremely randomized trees

```

Train  $K$  different trees with  $k$  bootstraps.
while desired depth do
    for number of tests do
        Pick  $m$  random variables
        Determine a random test
    end for
    Select test =  $\max\{InfoGain\}$ 
end while
Predict using ensemble

```

The smaller m is, the *more randomized* the trees are.

9.2 Boosting

Iterate training with focus on misclassified instances in the previous epoch.

With *AdaBoost*, the learner has a weight defined with:

3 Boosting

```
while error under threshold do
  for classifier do
    Train with more focus on misclassified instances
  end for
end while
Predict using ensemble
```

$$\text{Learner weight } \alpha_i = \frac{1}{2} \log\left(\frac{1 - \epsilon_i}{\epsilon_i}\right) \quad (72)$$

10 Neural Network

Generally can suffer from **overtraining** (*overfitting*) occurs when weights take on large magnitudes. Train with $P_{aram} = I \cdot O + b$ with respectively *Input*, *Output* sizes and *bias*.

- **Feed-forward**: output of layer j is input of $j + 1$.
Fully connected: all units in j are input of all units in $j + 1$.

$$\text{Sigmoid } \sigma(x) = \frac{e^x}{1 + e^x} \quad (73)$$

$$\text{Sigmoid derivate } \frac{\delta\sigma(x)}{\delta x} = \sigma(x)(1 - \sigma(x)) \quad (74)$$

$$\text{Hidden unit } h_i = \sigma(\bar{w}^T \bar{x} + b) \quad \forall i \quad (75)$$

4 FF-NN

```
 $\bar{h}^0 = 0$ 
for  $i = 1, \dots, H$  do
   $\bar{h}^i = \sigma(\bar{W}^i \bar{h}^{i-1} + \bar{b}^i)$ 
end for
 $y_{out} = \phi(\bar{W}^{out} \bar{h}^H + \bar{b}^{out})$ 
```

$$\text{Error} \quad J = \frac{1}{2}(\tilde{y} - y)^2 \quad (76)$$

$$\text{Error derivate (out)} \quad \frac{\partial J}{\partial \bar{w}_{out}} = \frac{\partial J}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial \bar{w}_{out}} = (\tilde{y} - y) \bar{h}_{last} = \delta_o \bar{h}_{last} \quad (77)$$

$$\text{Error derivate (hidden)} \quad \frac{\partial J}{\partial \bar{w}_j} = \frac{\partial J}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial \bar{w}_j} = \delta_o \frac{\partial \tilde{y}}{\partial \bar{w}_j} = \delta_o \frac{\partial \tilde{y}}{\partial \bar{h}_j} \frac{\partial \bar{h}_j}{\partial \bar{w}_j} \quad (78)$$

$$\Rightarrow \delta_o \bar{w}_{o,j} \sigma(\bar{w}^T \bar{x} + \bar{b})(1 - \sigma(\bar{w}^T \bar{x} + \bar{b})) \bar{x} \quad (79)$$

$$\Rightarrow \delta_{h,j} \bar{x} \quad (80)$$

5 Stochastic GDS

```

while no convergence do
  Pick a training  $\bar{x}$ 
  1 Feed  $\bar{x}$  and get  $y$ 
  2 Compute  $\frac{\partial J}{\partial \bar{w}_{out}}$ 
  for Hidden unit  $i$  do
    3 Compute share of correction  $\frac{\partial J}{\partial \bar{w}_j}$ 
  end for
  4 Update weights
    Hidden:  $\bar{w}_j = \bar{w}_j - \alpha \frac{\partial J}{\partial \bar{w}_j}$ 
    Output:  $\bar{w}_{out} = \bar{w}_{out} - \alpha \frac{\partial J}{\partial \bar{w}_{out}}$ 
end while

```

Claim 3. *Any function can be approximated to arbitrary accuracy by a network with 2 hidden layers.*

10.1 Activation function

Should be easy differentiable and **non-linear**.

- **Sigmoid, tanH:** easily saturate.
- **ReLU:** strong empirical results.
- **softplus:** smoother than ReLU but harder to train.
- **softmax:** generalizes multiclass tasks.

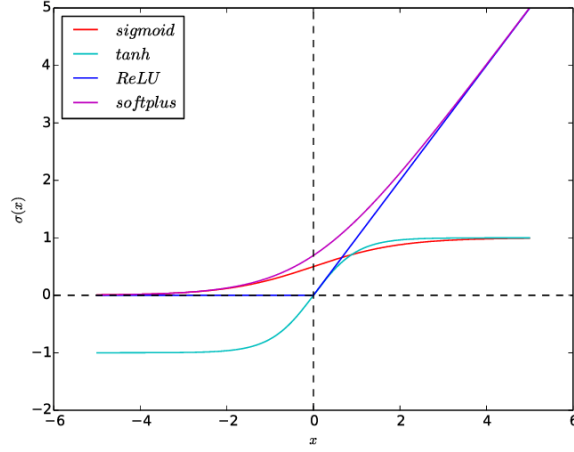


Figure 3: Most common activation functions

10.2 Backpropagation

Reverse-mode automatic differentiation (*RV-AD*) can efficiently compute the derivative of every node in a computation graph. **Very sensible to learning rates**. Empirically implemented in **AdamOpt**.

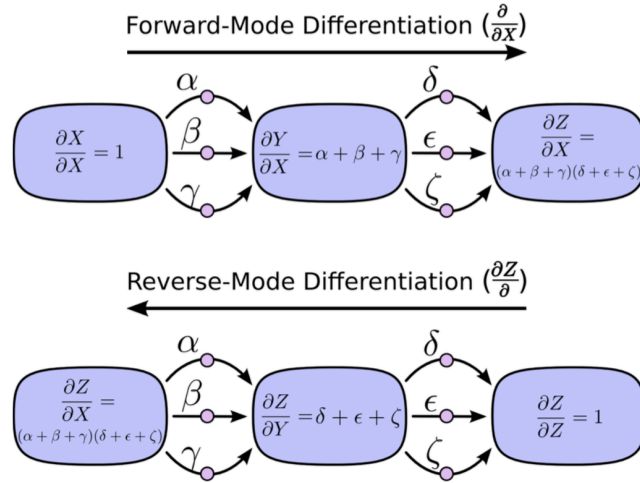


Figure 4: Differentiation paradigms from Hamilton (2019)

$$\text{Momentum} \quad M = \beta \Delta_{i-1} \bar{w} \quad (81)$$

$$\text{Momentum Update} \quad \Delta_i \bar{w} = \alpha \frac{\partial J}{\partial \bar{w}} + M \quad (82)$$

With momentum:

- **Pros:**
 - Avoid **small** local-minima.
 - Keep \bar{w} **moving** when error is flat.
- **Cons:**
 - Avoid **global minima**.
 - It's **an additional** parameter.

11 Convolutional Neural Network

Individual cortical neurons respond to stimuli only in a restricted region of the visual field known as the receptive field. The receptive fields of different neurons partially overlap such that they cover the entire visual field.

(Wikipedia, 2019)

Each layer transforms an input 3D **tensor** to an output 3D tensor using a differentiable function.

- **Very high-dimensional inputs**
- Multiple layers of inputs
- **Invariance:** light, rotations, translations

Main working scheme:

- **Local receptive field** for each first hidden unit (all channels), and **shared parameters** for each receptive field. For each conv: $(n \cdot m \cdot l)k$ parameters, with $n \times m$ filter, k output, and l input sizes.
- **Pooling:** aggregates result of convolutional layers.

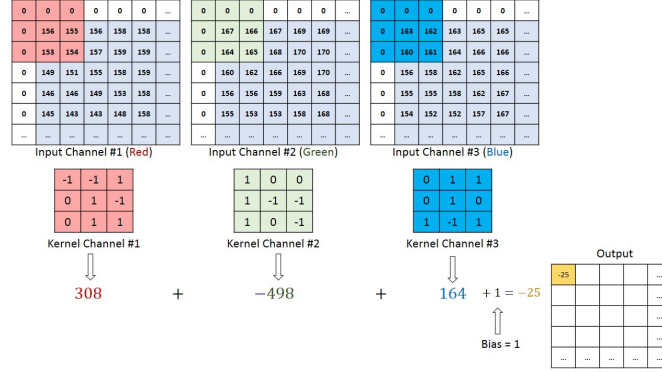


Figure 5: Convolution operation on a $M \times N \times 3$ image matrix with a $3 \times 3 \times 3$ Kernel. From TowardsDataScience (2019)

- **Stride**: spacing between receptive fields.

Definition 22. The **dropout** independently sets each hidden unit activity to zero with probability p

Definition 23. The **Batch Normalization** normalizes the input layer by adjusting and scaling the activations.

12 Recurrent Neural Network

RNN is a class of artificial neural network where connections between nodes form a directed graph along a **temporal sequence**. **Weights** are **shared** over time-steps. Output can be:

- **End of the sequence**: for instance, sentiment classification.
- **End of time-step**: for instance, generative language.

Recurrence can be based on:

- **Hidden states**: Elman RNN $\bar{h}_t = (\bar{W}\bar{h}_{t-1} + \bar{U}\bar{x}_t + \bar{b})$
- **Outputs**: Jordan RNN $\bar{h}_t = (\bar{W}\bar{o}_{t-1} + \bar{U}\bar{x}_t + \bar{b})$

$$\text{Output } \bar{o}_t = \phi(\bar{V}\bar{h}_t + \bar{c}) \quad (83)$$

Network are trained with **backpropagation** through time (*BPTT*). Sometimes the gradient is **truncated** for large sequences (see Figure 7).

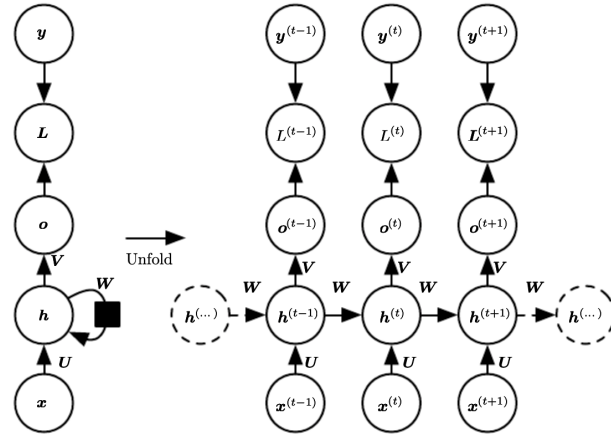


Figure 6: *RNN* scheme from Goodfellow et al. (2016)

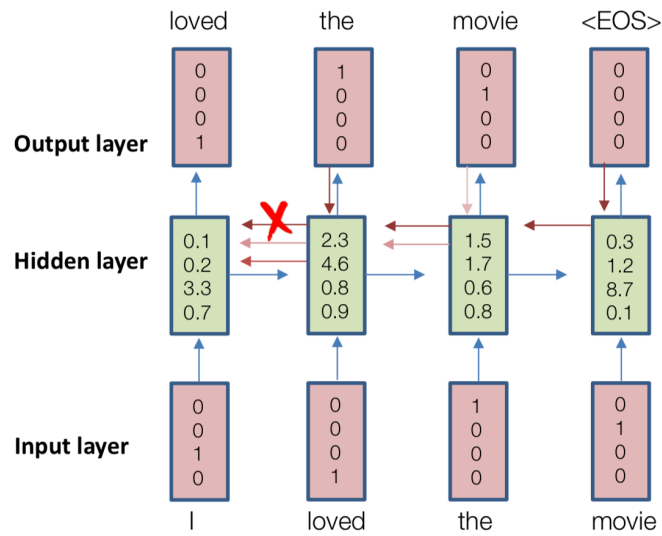


Figure 7: *RNN BPTT* from Hamilton (2019)

With **long term dependencies**: usually **exploding/vanishing** gradient:
 $\bar{h}_t = \bar{W} \bar{h}_{t-1}$ with $\bar{W} = \bar{Q} \bar{D} \bar{Q}^T$, and therefore $\bar{h}_t = \bar{Q} \bar{D}^d \bar{Q}^T \bar{h}_0$. Hence, use **Gradient Clipping**.

12.1 LSTM Units

Fix the **vanishing** gradient issue. Each cell is composed of **gates**, sigmoid layers to control information flow.

- **Hidden state tracking**: selects info from past for the *next* prediction.
- **Cell State**: selects info for future **predictions**.

Types of gates:

- **Forget gate**: amount of information to keep from previous cell.
- **Input gate**: what input is kept from previous cell.
- **Output gate**: what info from the cell is needed to predict.

12.2 Encoder-decoder RNN

Get different size output by **encoding** the input with a *RNN* - **context vector**-, and the output with a different *RNN*. The context vector can be enhanced with **attention**: namely, the decoder weights part of the context vector.

Definition 24. *The **teacher forcing** is a training procedure in which the model receives the ground truth output $y(t)$ as input at time $t+1$. (Goodfellow et al., 2016)*

References

- Bishop, C. M., 2006. Pattern Recognition and Machine Learning (Information Science and Statistics). Springer-Verlag, Berlin, Heidelberg.
- Goodfellow, I., Bengio, Y., Courville, A., Bengio, Y., 2016. Deep learning. Vol. 1. MIT press Cambridge.
- Hamilton, W., 2019. Comp551 lectures. University Lecture.
- TowardsDataScience, 2019. Towardsdatascience. Online.
- Wikipedia, 2019. Wikipedia.org. The Online Encyclopedia.