# Complete Mathematical Formulas in Machine Learning

## Table of Contents

---

## Foundational Mathematics

### Basic Notation

- **Scalar**: Single value, denoted as x, y, a, b

- **Vector**: Column of values, denoted as **x**, **y** (bold)

- **Matrix**: 2D array, denoted as **X**, **W** (bold capital)

- **Tensor**: n-dimensional array

### Summation and Product Notation

**Summation**:

```
Σᵢ₌₁ⁿ xᵢ = x₁ + x₂ + ... + xₙ
```

$$\sum_{i=1}^{n} x_i = x_1 + x_2 + \ldots + x_n$$

**Product**:

$$\prod_{i=1}^{n} x_i = x_1 \times x_2 \times \ldots \times x_n$$

---

# Statistical Measures

## Mean (Average)

**Arithmetic Mean**:

$$\mu = \bar{x} = (1/n) \times \sum_{i=1}^{n} x_i$$

- $\mu$ (mu) or $\bar{x}$ represents the mean
- $n$ is the number of observations
- $x_i$ is the i-th observation

## Variance

**Population Variance**:

$$\sigma^2 = (1/n) \times \sum_{i=1}^{n} (x_i - \mu)^2$$

**Sample Variance** (Bessel's correction):

$$s^2 = (1/(n-1)) \times \sum_{i=1}^{n} (x_i - \bar{x})^2$$

## Standard Deviation

```
σ = √(σ²)  (population)
s = √(s²)  (sample)
```

## Covariance

**Between two variables X and Y**:

$$Cov(X,Y) = (1/n) \times \sum_{i=1}^{n} (x_i - \mu_x)(y_i - \mu_Y)$$

## Correlation Coefficient

**Pearson Correlation**:

$$\rho(X,Y) = Cov(X,Y) / (\sigma_x \times \sigma_Y)$$

- Range: [-1, 1]

- -1: perfect negative correlation

- 0: no linear correlation

- 1: perfect positive correlation

## Standardization (Z-score)

```
z = (x - μ) / σ
```

- Transforms data to have $\mu = 0$ and $\sigma = 1$

---

# Linear Algebra in ML

## Vector Operations

### Dot Product (Inner Product):

```
a · b = Σᵢ₌₁ⁿ aᵢbᵢ = a₁b₁ + a₂b₂ + ... + aₙbₙ
```

### Vector Norm (Length):

- L1 Norm (Manhattan): $\|x\|_1 = \Sigma_i |x_i|$
- L2 Norm (Euclidean): $\|x\|_2 = \sqrt{\Sigma_i x_i^2}$
- L∞ Norm (Max): $\|x\|_\infty = \max_i |x_i|$

## Matrix Operations

### Matrix Multiplication:

```
C = AB where Cᵢⱼ = Σₖ AᵢₖBₖⱼ
```

### Matrix Transpose:

```
(Aᵀ)ᵢⱼ = Aⱼᵢ
```

### Matrix Inverse:

```
AA⁻¹ = A⁻¹A = I
```

where I is the identity matrix

## Eigenvalues and Eigenvectors

For matrix A and vector v:

```
Av = λv
```

- λ is the eigenvalue
- v is the eigenvector

**Characteristic Equation**:

```
det(A - λI) = 0
```

---

# Calculus in ML

## Derivatives

**Basic Rules**:

- Power Rule: $d/dx(x^n) = nx^{n-1}$
- Chain Rule: $d/dx[f(g(x))] = f'(g(x)) \times g'(x)$
- Product Rule: $d/dx[f(x)g(x)] = f'(x)g(x) + f(x)g'(x)$

## Partial Derivatives

For function f(x,y):

```
∂f/∂x = partial derivative with respect to x
∂f/∂y = partial derivative with respect to y
```

## Gradient

For function $f(x_1, x_2, ..., x_n)$:

```
∇f = [∂f/∂x₁, ∂f/∂x₂, ..., ∂f/∂xₙ]ᵀ
```

## Gradient Descent Update Rule

```
θₜ₊₁ = θₜ - α∇f(θₜ)
```

- θ: parameters
- α: learning rate
- ∇f: gradient of loss function

## Hessian Matrix

Second-order partial derivatives:

```
H(f)ᵢⱼ = ∂²f/(∂xᵢ∂xⱼ)
```

---

# Probability Theory

## Basic Probability

```
P(A) ∈ [0, 1]
P(Ω) = 1  (probability of sample space)
P(∅) = 0  (probability of empty set)
```

## Conditional Probability

```
P(A|B) = P(A ∩ B) / P(B)
```

- P(A|B): probability of A given B

## Bayes' Theorem

```
P(A|B) = [P(B|A) × P(A)] / P(B)
```

**Extended form**:

```
P(A|B) = [P(B|A) × P(A)] / [Σᵢ P(B|Aᵢ) × P(Aᵢ)]
```

## Probability Distributions

### Bernoulli Distribution:

```
P(X = x) = pˣ(1-p)¹⁻ˣ for x ∈ {0,1}
```

### Binomial Distribution:

```
P(X = k) = C(n,k) × pᵏ × (1-p)ⁿ⁻ᵏ
```

where C(n,k) = n!/(k!(n-k)!)

### Normal (Gaussian) Distribution:

```
f(x) = (1/√(2πσ²)) × exp(-(x-μ)²/(2σ²))
```

**Multivariate Normal**:

$$f(x) = (1/((2\pi)^{(k/2)}|\Sigma|^{(1/2)})) \times \exp(-\tfrac{1}{2}(x-\mu)^T\Sigma^{-1}(x-\mu))$$

- $\Sigma$: covariance matrix
- $|\Sigma|$: determinant of $\Sigma$

## Expectation and Variance

**Expectation** (Expected Value):

- Discrete: $E[X] = \Sigma_x\, x \times P(X = x)$
- Continuous: $E[X] = \int x \times f(x)dx$

**Variance**:

$$Var(X) = E[(X - E[X])^2] = E[X^2] - (E[X])^2$$

---

# Information Theory

## Entropy

**For discrete variable X**:

$$H(X) = -\Sigma_i\, P(x_i) \times \log_2(P(x_i))$$

- Measures uncertainty/information content
- Units: bits ($\log_2$) or nats (ln)

## Cross-Entropy

**Between distributions P and Q**:

$$H(P,Q) = -\Sigma_i\, P(x_i) \times \log(Q(x_i))$$

## Kullback-Leibler (KL) Divergence

$$KL(P||Q) = \Sigma_i\, P(x_i) \times \log(P(x_i)/Q(x_i))$$

- Measures difference between distributions
- $KL(P||Q) \neq KL(Q||P)$ (not symmetric)

## Mutual Information

```
I(X;Y) = Σ_x Σ_y P(x,y) × log(P(x,y)/(P(x)P(y)))
```

---

## Loss Functions

### Regression Loss Functions

**Mean Squared Error (MSE)**:

```
MSE = (1/n) × Σ_{i=1}^n (y_i - ŷ_i)²
```

**Mean Absolute Error (MAE)**:

```
MAE = (1/n) × Σ_{i=1}^n |y_i - ŷ_i|
```

**Huber Loss** (Robust to outliers):

```
L_δ(y,ŷ) = {
    ½(y-ŷ)²           if |y-ŷ| ≤ δ
    δ|y-ŷ| - ½δ²      if |y-ŷ| > δ
}
```

### Classification Loss Functions

**Binary Cross-Entropy** (Log Loss):

```
BCE = -(1/n) × Σ_{i=1}^n [y_i log(ŷ_i) + (1-y_i)log(1-ŷ_i)]
```

**Categorical Cross-Entropy** (Multi-class):

```
CCE = -(1/n) × Σ_{i=1}^n Σ_{j=1}^m y_{ij} log(ŷ_{ij})
```

- m: number of classes
- $y_{ij}$: 1 if sample i belongs to class j, 0 otherwise

**Hinge Loss** (SVM):

```
L = max(0, 1 - y × ŷ)
```

---

## Optimization Algorithms

### Gradient Descent Variants

**Batch Gradient Descent**:

```
θ = θ - α × (1/n) × Σᵢ₌₁ⁿ ∇θL(xᵢ, yᵢ, θ)
```

**Stochastic Gradient Descent (SGD)**:

```
θ = θ - α × ∇θL(xᵢ, yᵢ, θ)
```

**Mini-batch Gradient Descent**:

```
θ = θ - α × (1/m) × Σᵢ₌₁ᵐ ∇θL(xᵢ, yᵢ, θ)
```

## Advanced Optimizers

**Momentum**:

$$v_t = \beta v_{t-1} + \alpha \nabla_\theta L$$
$$\theta_t = \theta_{t-1} - v_t$$

**RMSprop**:

$$s_t = \beta s_{t-1} + (1-\beta)(\nabla_\theta L)^2$$
$$\theta_t = \theta_{t-1} - \alpha \times \nabla_\theta L/\sqrt{(s_t + \varepsilon)}$$

**Adam** (Adaptive Moment Estimation):

$$m_t = \beta_1 m_{t-1} + (1-\beta_1)\nabla_\theta L \qquad \text{(1st moment)}$$
$$v_t = \beta_2 v_{t-1} + (1-\beta_2)(\nabla_\theta L)^2 \qquad \text{(2nd moment)}$$
$$\hat{m}_t = m_t/(1-\beta_1^t) \qquad \text{(bias correction)}$$
$$\hat{v}_t = v_t/(1-\beta_2^t) \qquad \text{(bias correction)}$$
$$\theta_t = \theta_{t-1} - \alpha \times \hat{m}_t/(\sqrt{\hat{v}_t} + \varepsilon)$$

---

# Linear Models Mathematics

## Linear Regression

**Model**:

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_n x_n = x^T \beta$$

**Normal Equation** (Closed-form solution):

$$\beta = (X^TX)^{-1}X^Ty$$

**Cost Function** (MSE):

$$J(\beta) = (1/2n) \times \Sigma_{i=1}^{n} (y_i - x_i^T\beta)^2$$

**Gradient**:

$$\nabla\beta J = -(1/n) \times X^T(y - X\beta)$$

## Ridge Regression (L2 Regularization)

**Cost Function**:

$$J(\beta) = (1/2n) \times \Sigma_{i=1}^{n} (y_i - x_i^T\beta)^2 + \lambda||\beta||_2^2$$

**Solution**:

$$\beta = (X^TX + \lambda I)^{-1}X^Ty$$

## Lasso Regression (L1 Regularization)

**Cost Function**:

$$J(\beta) = (1/2n) \times \Sigma_{i=1}^{n} (y_i - x_i^T\beta)^2 + \lambda||\beta||_1$$

- No closed-form solution
- Solved using coordinate descent or proximal gradient

## Logistic Regression

**Sigmoid Function**:

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

**Model**:

$$P(y=1|x) = \sigma(x^T\beta) = 1/(1 + e^{-x^T\beta})$$

**Log-Likelihood**:

```
LL = Σᵢ₌₁ⁿ [yᵢlog(σ(xᵢᵀβ)) + (1-yᵢ)log(1-σ(xᵢᵀβ))]
```

**Gradient**:

```
∇βLL = Xᵀ(y - σ(Xβ))
```

---

# Tree-Based Models Mathematics

## Decision Trees

### Gini Impurity:

```
Gini = 1 - Σⱼ₌₁ᶜ pⱼ²
```

- c: number of classes
- $p_j$: proportion of samples in class j

### Entropy:

```
Entropy = -Σⱼ₌₁ᶜ pⱼlog₂(pⱼ)
```

### Information Gain:

```
IG = Entropy(parent) - Σᵢ (nᵢ/n) × Entropy(childᵢ)
```

- $n_i$: number of samples in child i
- n: total samples in parent

### Variance Reduction (for regression):

```
VR = Var(parent) - Σᵢ (nᵢ/n) × Var(childᵢ)
```

## Random Forest

### Prediction (Regression):

```
ŷ = (1/B) × Σᵦ₌₁ᴮ Tᵦ(x)
```

- B: number of trees
- $T_\beta$: prediction from tree b

**Prediction** (Classification):

```
ŷ = mode{T₁(x), T₂(x), ..., Tᴮ(x)}
```

**Feature Importance**:

```
Importanceⱼ = (1/B) × Σβ=₁ᴮ Σₜ∈β 1(v(t)=j) × p(t)Δᵢₜ
```

- v(t): variable used at node t
- p(t): proportion of samples reaching node t
- $\Delta_{it}$: impurity decrease at node t

## Gradient Boosting

### Additive Model:

```
F(x) = Σm=₁ᴹ γₘhₘ(x)
```

### Update Rule:

```
Fₘ(x) = Fₘ₋₁(x) + γₘhₘ(x)
```

### Gradient Calculation:

```
rᵢₘ = -[∂L(yᵢ, F(xᵢ))/∂F(xᵢ)]_{F=Fₘ₋₁}
```

---

# Support Vector Machines Mathematics

## Hard Margin SVM

### Objective:

```
minimize: ½||w||²
subject to: yᵢ(wᵀxᵢ + b) ≥ 1 for all i
```

## Soft Margin SVM

### Objective:

```
minimize: ½||w||² + C × Σᵢ ξᵢ
subject to: yᵢ(wᵀxᵢ + b) ≥ 1 - ξᵢ
            ξᵢ ≥ 0
```

- $\xi_i$: slack variables
- C: regularization parameter

## Kernel Trick

**Kernel Function**:

```
K(x, x') = φ(x)ᵀφ(x')
```

**Common Kernels**:

- Linear: $K(x, x') = x^T x'$
- Polynomial: $K(x, x') = (\gamma x^T x' + r)^d$
- RBF (Gaussian): $K(x, x') = \exp(-\gamma||x - x'||^2)$
- Sigmoid: $K(x, x') = \tanh(\gamma x^T x' + r)$

**Dual Form Prediction**:

```
f(x) = Σᵢ₌₁ⁿ αᵢyᵢK(xᵢ, x) + b
```

---

# Neural Network Mathematics

## Forward Propagation

**Linear Transformation**:

```
z⁽¹⁾ = W⁽¹⁾a⁽¹⁻¹⁾ + b⁽¹⁾
```

**Activation**:

```
a⁽¹⁾ = g⁽¹⁾(z⁽¹⁾)
```

## Activation Functions

**ReLU** (Rectified Linear Unit):

```
ReLU(z) = max(0, z)
```

**Leaky ReLU**:

```
LeakyReLU(z) = max(αz, z)  where α ≈ 0.01
```

**Sigmoid**:

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

**Tanh**:

$$\tanh(z) = (e^z - e^{-z})/(e^z + e^{-z})$$

**Softmax** (Multi-class output):

$$\text{softmax}(z_i) = e^{(z_i)} / \Sigma_j \, e^{(z_j)}$$

# Backpropagation

## Chain Rule Application:

$$\partial L/\partial W^{(l)} = \partial L/\partial z^{(l)} \times \partial z^{(l)}/\partial W^{(l)}$$

## Error Propagation:

$$\delta^{(l)} = (W^{(l+1)})^T \delta^{(l+1)} \odot g'(z^{(l)})$$

- $\odot$: element-wise multiplication

## Weight Update:

$$W^{(l)} = W^{(l)} - \alpha \times \delta^{(l)}(a^{(l-1)})^T$$

# Batch Normalization

## Normalization:

$$\hat{x}_i = (x_i - \mu B)/\sqrt{(\sigma B^2 + \varepsilon)}$$

## Scale and Shift:

$$y_i = \gamma \hat{x}_i + \beta$$

- γ, β: learnable parameters

## Dropout

**Training**:

```
a⁽¹⁾ = a⁽¹⁾ ⊙ m⁽¹⁾ / p
```

- $m^{(i)}$: binary mask (Bernoulli(p))
- p: keep probability

---

# Clustering Mathematics

## K-Means

**Objective Function**:

```
J = Σᵢ₌₁ⁿ Σₖ₌₁ᴷ rᵢₖ||xᵢ - μₖ||²
```

- $r_{ik}$: 1 if $x_i$ belongs to cluster k, 0 otherwise

**Update Rules**:

```
μₖ = (Σᵢ rᵢₖxᵢ)/(Σᵢ rᵢₖ)  (centroid update)
rᵢₖ = 1 if k = argmin_j ||xᵢ - μⱼ||²  (assignment)
```

## DBSCAN

**Core Point**:

```
|Nε(p)| ≥ minPts
```

- Nε(p) = {q ∈ D | dist(p,q) ≤ ε}

**Density-Reachable**:

- p is reachable from q if there's a chain of core points

## Gaussian Mixture Models

**Probability Model**:

```
p(x) = Σₖ₌₁ᴷ πₖ × N(x|μₖ, Σₖ)
```

**E-step** (Expectation):

$$\gamma_{ik} = (\pi_k \times N(x_i|\mu_k, \Sigma_k))/(\Sigma_j \pi_j \times N(x_i|\mu_j, \Sigma_j))$$

**M-step** (Maximization):

$$\pi_k = (1/n) \times \Sigma_i \gamma_{ik}$$
$$\mu_k = (\Sigma_i \gamma_{ik}x_i)/(\Sigma_i \gamma_{ik})$$
$$\Sigma_k = (\Sigma_i \gamma_{ik}(x_i - \mu_k)(x_i - \mu_k)^\top)/(\Sigma_i \gamma_{ik})$$

---

# Dimensionality Reduction Mathematics

## Principal Component Analysis (PCA)

**Objective**: Find projection that maximizes variance

**Covariance Matrix**:

$$C = (1/n) \times X^\top X$$

**Eigendecomposition**:

$$C = V\Lambda V^\top$$

- V: eigenvectors (principal components)
- Λ: eigenvalues (diagonal matrix)

**Projection**:

$$Z = XV_k$$

- $V_k$: first k eigenvectors

**Reconstruction**:

$$\hat{X} = ZV_k{}^\top$$

**Explained Variance Ratio**:

$$EVR = \lambda_i / \Sigma_j \lambda_j$$

## Linear Discriminant Analysis (LDA)

**Between-class Scatter**:

$$S_B = \Sigma_k\, n_k(\mu_k - \mu)(\mu_k - \mu)^\top$$

**Within-class Scatter**:

$$S_W = \Sigma_k\, \Sigma_{i \in C_k}\, (x_i - \mu_k)(x_i - \mu_k)^\top$$

**Objective**:

```
maximize: (wᵀSʙw)/(wᵀSww)
```

**Solution**: Eigenvectors of $S_W^{-1}S_B$

## t-SNE

**Joint Probability (High-dimensional)**:

```
pᵢⱼ = (pⱼ|ᵢ + pᵢ|ⱼ)/(2n)
```

where:

```
pⱼ|ᵢ = exp(-||xᵢ - xⱼ||²/(2σᵢ²)) / Σₖ≠ᵢ exp(-||xᵢ - xₖ||²/(2σᵢ²))
```

**Joint Probability (Low-dimensional)**:

```
qᵢⱼ = (1 + ||yᵢ - yⱼ||²)⁻¹ / Σₖ≠ₗ (1 + ||yₖ - yₗ||²)⁻¹
```

**Cost Function** (KL divergence):

```
C = KL(P||Q) = Σᵢ Σⱼ pᵢⱼ log(pᵢⱼ/qᵢⱼ)
```

---

# Evaluation Metrics Mathematics

## Classification Metrics

### Confusion Matrix Elements:

- TP: True Positives
- TN: True Negatives
- FP: False Positives
- FN: False Negatives

**Accuracy**:

```
Accuracy = (TP + TN)/(TP + TN + FP + FN)
```

**Precision**:

```
Precision = TP/(TP + FP)
```

**Recall** (Sensitivity, True Positive Rate):

```
Recall = TP/(TP + FN)
```

**Specificity** (True Negative Rate):

```
Specificity = TN/(TN + FP)
```

**F1 Score**:

```
F1 = 2 × (Precision × Recall)/(Precision + Recall)
```

**F-beta Score**:

```
Fβ = (1 + β²) × (Precision × Recall)/((β² × Precision) + Recall)
```

**Matthews Correlation Coefficient**:

```
MCC = (TP×TN - FP×FN)/√((TP+FP)(TP+FN)(TN+FP)(TN+FN))
```

## ROC and AUC

**ROC Curve**: Plot of TPR vs FPR at various thresholds

**AUC** (Area Under ROC Curve):

```
AUC = ∫₀¹ TPR(FPR) dFPR
```

**Approximation** (Trapezoidal Rule):

```
AUC ≈ Σᵢ ½(TPRᵢ + TPRᵢ₊₁)(FPRᵢ₊₁ - FPRᵢ)
```

## Regression Metrics

**R² Score** (Coefficient of Determination):

```
R² = 1 - (SS_res/SS_tot)
```

where:

```
SS_res = Σᵢ (yᵢ - ŷᵢ)²  (residual sum of squares)
SS_tot = Σᵢ (yᵢ - ȳ)²   (total sum of squares)
```

**Adjusted R²**:

```
R²_adj = 1 - [(1-R²)(n-1)/(n-p-1)]
```

- p: number of predictors

**Mean Absolute Percentage Error (MAPE)**:

```
MAPE = (100/n) × Σᵢ |yᵢ - ŷᵢ|/|yᵢ|
```

## Clustering Metrics

### Silhouette Score:

```
s(i) = (b(i) - a(i))/max(a(i), b(i))
```

where:

- $a(i)$: average distance to points in same cluster
- $b(i)$: minimum average distance to points in other clusters

### Davies-Bouldin Index:

```
DB = (1/k) × Σᵢ₌₁ᵏ maxⱼ≠ᵢ [(σᵢ + σⱼ)/d(cᵢ, cⱼ)]
```

- $\sigma_i$: average distance of points in cluster i to centroid
- $d(c_i, c_j)$: distance between centroids

### Calinski-Harabasz Index:

```
CH = [tr(Bₖ)/(k-1)] / [tr(Wₖ)/(n-k)]
```

- $B_k$: between-group dispersion matrix

- $W_k$: within-group dispersion matrix

---

# Additional Important Formulas

## Regularization

**Elastic Net** (L1 + L2):

```
J(θ) = L(θ) + λ₁||θ||₁ + λ₂||θ||₂²
```

## Distance Metrics

**Euclidean Distance**:

```
d(x,y) = √(Σᵢ (xᵢ - yᵢ)²)
```

**Manhattan Distance**:

```
d(x,y) = Σᵢ |xᵢ - yᵢ|
```

**Cosine Similarity**:

```
cos(x,y) = (x·y)/(||x||₂ × ||y||₂)
```

**Minkowski Distance**:

```
d(x,y) = (Σᵢ |xᵢ - yᵢ|ᵖ)^(1/p)
```

## Bias-Variance Decomposition

**Total Error**:

```
E[(y - ŷ)²] = Bias² + Variance + Irreducible Error
```

where:

```
Bias = E[ŷ] - y
Variance = E[(ŷ - E[ŷ])²]
```

## Cross-Validation Error

**k-Fold CV Error**:

```
CV(k) = (1/k) × Σ_{i=1}^{k} L(h_i, D_i)
```

- $h_i$: model trained on all folds except i
- $D_i$: validation data from fold i

---

## Time Series Analysis

### Autoregressive (AR) Model

$$y_t = c + \Sigma_{i=1}^{p} \phi_i y_{t-i} + \varepsilon_t$$

- $\varphi_i$: AR coefficients
- p: order of AR model

### Moving Average (MA) Model

$$y_t = \mu + \varepsilon_t + \Sigma_{i=1}^{g} \theta_i \varepsilon_{t-i}$$

- $\theta_i$: MA coefficients
- q: order of MA model

### ARIMA Model

$$(1 - \Sigma_{i=1}^{p} \phi_i L^i)(1-L)^d y_t = (1 + \Sigma_{i=1}^{g} \theta_i L^i)\varepsilon_t$$

- L: lag operator
- d: degree of differencing

### Exponential Smoothing

**Simple**: $\alpha y_t + (1-\alpha)\hat{y}_{t-1}$ **Double**: $\alpha y_t + (1-\alpha)(\hat{y}_{t-1} + \hat{b}_{t-1})$ **Triple (Holt-Winters)**: Includes seasonal component

### Autocorrelation Function (ACF)

$$\rho_k = Cov(y_t, y_{t-k}) / Var(y_t)$$

### Partial Autocorrelation Function (PACF)

Correlation between $y_t$ and $y_{t-k}$ after removing effects of intermediate lags

---

## Bayesian Methods

## Bayes' Rule (Full Form)

$$P(\theta|D) = P(D|\theta)P(\theta) / P(D)$$

- $P(\theta|D)$: Posterior
- $P(D|\theta)$: Likelihood
- $P(\theta)$: Prior
- $P(D)$: Evidence

## Maximum A Posteriori (MAP)

$$\theta_{MAP} = \text{argmax}_\theta\ P(\theta|D) = \text{argmax}_\theta\ P(D|\theta)P(\theta)$$

## Bayesian Linear Regression

**Posterior**:

$$P(w|D) = N(w|\mu_n,\ \Sigma_n)$$

where:

$$\Sigma_n = (\Sigma_0^{-1} + \beta X^\top X)^{-1}$$
$$\mu_n = \Sigma_n(\Sigma_0^{-1}\mu_0 + \beta X^\top y)$$

## Variational Inference

**ELBO (Evidence Lower Bound)**:

$$L(q) = E_q[\log P(X,Z)] - E_q[\log q(Z)]$$

## Markov Chain Monte Carlo (MCMC)

**Metropolis-Hastings Accept Probability**:

$$\alpha = \min(1,\ P(\theta'|D)Q(\theta|\theta') / P(\theta|D)Q(\theta'|\theta))$$

---

# Convolutional Neural Networks

## Convolution Operation

$$(f * g)[n] = \sum_m f[m] \times g[n-m]$$

## 2D Convolution (Images)

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

- I: input image
- K: kernel/filter

## Pooling

**Max Pooling**:

$$y = \max(x_1, x_2, ..., x_n) \text{ in pooling window}$$

**Average Pooling**:

$$y = (1/n)\sum_i x_i \text{ in pooling window}$$

## Output Size Calculation

$$\text{Output Size} = \lfloor (W - F + 2P)/S \rfloor + 1$$

- W: input size
- F: filter size
- P: padding
- S: stride

## Number of Parameters

**Conv Layer**: $(F \times F \times C_{in} + 1) \times C_{out}$ **FC Layer**: $(N_{in} + 1) \times N_{out}$

---

# Recurrent Neural Networks

## Vanilla RNN

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$
$$y_t = W_{hy}h_t + b_y$$

## LSTM (Long Short-Term Memory)

**Forget Gate**:

```
f_t = σ(Wf·[h_{t-1}, x_t] + bf)
```

**Input Gate**:

```
i_t = σ(Wi·[h_{t-1}, x_t] + bi)
C̃_t = tanh(WC·[h_{t-1}, x_t] + bC)
```

**Cell State Update**:

```
C_t = f_t * C_{t-1} + i_t * C̃_t
```

**Output Gate**:

```
o_t = σ(Wo·[h_{t-1}, x_t] + bo)
h_t = o_t * tanh(C_t)
```

## GRU (Gated Recurrent Unit)

**Update Gate**: $z_t = σ(Wz·[h_{t-1}, x_t])$ **Reset Gate**: $r_t = σ(Wr·[h_{t-1}, x_t])$ **Hidden State**: $h_t = (1-z_t)h_{t-1} + z_t\tilde{h}_t$

---

# Transformer Architecture

## Scaled Dot-Product Attention

```
Attention(Q,K,V) = softmax(QKᵀ/√d_k)V
```

- Q: Query matrix
- K: Key matrix
- V: Value matrix
- $d_k$: dimension of keys

## Multi-Head Attention

```
MultiHead(Q,K,V) = Concat(head_1,...,head_h)Wᴼ
```

where:

```
head_i = Attention(QW_iᵠ, KW_iᴷ, VW_iⱽ)
```

## Positional Encoding

```
PE(pos,2i) = sin(pos/10000^(2i/d_model))
PE(pos,2i+1) = cos(pos/10000^(2i/d_model))
```

## Layer Normalization

```
LN(x) = γ × (x-μ)/√(σ²+ε) + β
```

---

# Recommender Systems

## Collaborative Filtering

**Matrix Factorization**:

```
R ≈ PQᵀ
```

- R: user-item rating matrix
- P: user feature matrix
- Q: item feature matrix

**Objective Function**:

```
min Σ(i,j)∈K (rᵢⱼ - pᵢᵀqⱼ)² + λ(||pᵢ||² + ||qⱼ||²)
```

## Cosine Similarity

```
sim(u,v) = (Σᵢ rᵤᵢrᵥᵢ) / (√Σᵢ rᵤᵢ² × √Σᵢ rᵥᵢ²)
```

## Pearson Correlation

```
sim(u,v) = Σᵢ(rᵤᵢ-r̄ᵤ)(rᵥᵢ-r̄ᵥ) / √(Σᵢ(rᵤᵢ-r̄ᵤ)² × Σᵢ(rᵥᵢ-r̄ᵥ)²)
```

---

# Reinforcement Learning Mathematics

## Bellman Equation

**For Value Function**:

```
V(s) = max_a [R(s,a) + γΣₛ' P(s'|s,a)V(s')]
```

**For Q-Function**:
```

$$Q(s,a) = R(s,a) + \gamma\Sigma_{s'} P(s'|s,a)\max_{a'} Q(s',a')$$

## Policy Gradient Theorem

$$\nabla\theta J(\theta) = E[\nabla\theta \log \pi\theta(a|s) \times Q\pi(s,a)]$$

## Temporal Difference Learning

### TD(0):

$$V(s_t) \leftarrow V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

### TD(λ):

$$V(s) \leftarrow V(s) + \alpha\delta_t e_t(s)$$

where eligibility trace:

$$e_t(s) = \gamma\lambda e_{t-1}(s) + 1(s=s_t)$$

## Actor-Critic

**Critic Update**: Update value function V(s) **Actor Update**: $\nabla\theta J = E[\nabla\theta \log \pi\theta(a|s) \times A(s,a)]$ where advantage A(s,a) = Q(s,a) - V(s)

---

# Statistical Hypothesis Testing

## Z-Test

$$z = (\bar{x} - \mu) / (\sigma/\sqrt{n})$$

## T-Test

**One Sample**:

$$t = (\bar{x} - \mu) / (s/\sqrt{n})$$

**Two Sample** (Equal Variance):

$$t = (\bar{x}_1 - \bar{x}_2) / (s_p\sqrt{1/n_1 + 1/n_2})$$

where pooled variance:

$$s_p^2 = ((n_1-1)s_1^2 + (n_2-1)s_2^2) / (n_1+n_2-2)$$

## Chi-Square Test

$$\chi^2 = \Sigma_i (O_i - E_i)^2 / E_i$$

- $O_i$: observed frequency
- $E_i$: expected frequency

## ANOVA F-Statistic

$$F = MSB / MSW = (SSB/(k-1)) / (SSW/(n-k))$$

- MSB: mean square between groups
- MSW: mean square within groups

## p-value

```
p-value = P(|Test Statistic| ≥ |Observed Value| | H₀)
```

---

# Ensemble Methods Mathematics

## Bagging

**Bootstrap Sampling**: Sample n instances with replacement **Aggregation**: $\hat{y} = (1/B)\Sigma_B f_\beta(x)$

## AdaBoost

**Sample Weight Update**:

$$w_i^{(t+1)} = w_i^{(t)} \times \exp(-\alpha_t y_i h_t(x_i))$$

**Classifier Weight**:

$$\alpha_t = \tfrac{1}{2} \log((1-\varepsilon_t)/\varepsilon_t)$$

**Final Prediction**:

$$H(x) = \text{sign}(\Sigma_t \alpha_t h_t(x))$$

## Gradient Boosting Mathematics

**Pseudo-Residuals**:

$$r_{im} = -[\partial L(y_i, f(x_i))/\partial f(x_i)]_{\{f=f_{m-1}\}}$$

**Line Search**:

$$\gamma_m = \text{argmin}\_\gamma \; \Sigma_i \; L(y_i, \; f_{m-1}(x_i) + \gamma h_m(x_i))$$

## XGBoost Objective

$$L = \Sigma_i \; l(y_i, \hat{y}_i) + \Sigma_k \; \Omega(f_k)$$

where:

$$\Omega(f) = \gamma T + \tfrac{1}{2}\lambda ||w||^2$$

- T: number of leaves
- w: leaf weights

---

# Graph Neural Networks

## Graph Convolution

$$H^{(l+1)} = \sigma(\tilde{D}^{-1/2}\tilde{A}\tilde{D}^{-1/2}H^{(l)}W^{(l)})$$

- $\tilde{A} = A + I$ (adjacency matrix with self-loops)
- $\tilde{D}$: degree matrix of $\tilde{A}$

## Message Passing

$$h_i^{(k+1)} = \sigma(W_{self}h_i^{(k)} + \Sigma_j \in N(i) \; W_{neigh}h_j^{(k)})$$

## Graph Attention

$$\alpha_{ij} = \text{softmax}_j(\text{LeakyReLU}(a^\top[Wh_i || Wh_j]))$$

---

# Advanced Optimization

## Newton's Method

$$\theta_{t+1} = \theta_t - H^{-1}\nabla f(\theta_t)$$

- H: Hessian matrix

## L-BFGS

Approximates inverse Hessian using limited memory

## Conjugate Gradient

$$d_{k+1} = -g_{k+1} + \beta_k d_k$$

where:

$$\beta_k = g_{k+1}^T g_{k+1} / g_k^T g_k \text{ (Fletcher-Reeves)}$$

## Natural Gradient

$$\theta_{t+1} = \theta_t - \alpha F^{-1}\nabla L(\theta_t)$$

- F: Fisher Information Matrix

---

# Sampling Methods

## Rejection Sampling

Accept x with probability:

```
p(accept) = f(x) / (M×g(x))
```

## Importance Sampling

```
E_p[f(x)] = E_q[f(x)p(x)/q(x)]
```

## Gibbs Sampling

Sample each variable conditioned on others:

$$x_i^{(t+1)} \sim P(x_i | x_1^{(t+1)},...,x_{i-1}^{(t+1)}, x_{i+1}^{(t)},...,x_n^{(t)})$$

---

# Online Learning

## Stochastic Gradient Descent

$$\theta_{t+1} = \theta_t - \eta_t \nabla l(x_t, y_t, \theta_t)$$

## Online Gradient Descent Regret Bound

$$\text{Regret} \leq ||\theta^*||^2/(2\eta) + \eta \times \Sigma_t ||g_t||^2/2$$

## Exponential Weighted Average

$$w_{t+1,i} = w_{t,i} \times \exp(-\eta l_{t,i}) / Z_t$$

---

# Semi-Supervised Learning

## Self-Training Loss

$$L = \Sigma_l \, L(x_l, y_l) + \lambda \Sigma_u \, L(x_u, \hat{y}_u)$$

## Graph-Based SSL

**Label Propagation**:

$$f_{t+1} = \alpha W f_t + (1-\alpha)y$$

## Co-Training

Train two classifiers on different views:

$$h_1: X_1 \rightarrow Y$$
$$h_2: X_2 \rightarrow Y$$

---

# Probabilistic Graphical Models

## Hidden Markov Models (HMM)

**Forward Algorithm**:

$$\alpha_t(j) = [\Sigma_i \, \alpha_{t-1}(i) a_{ij}] \times b_j(o_t)$$

- $a_{ij}$: transition probability
- $b_j(o_t)$: emission probability

**Backward Algorithm**:

$$\beta_t(i) = \Sigma_j\ a_{ij}b_j(o_{t+1})\beta_{t+1}(j)$$

**Viterbi Algorithm**:

$$\delta_t(j) = \max\_i\ [\delta_{t-1}(i) \times a_{ij}] \times b_j(o_t)$$

**Baum-Welch Update**:

$$a_{ij} = \Sigma_t\ \xi_t(i,j)\ /\ \Sigma_t\ \gamma_t(i)$$

## Conditional Random Fields (CRF)

$$P(y|x) = (1/Z(x)) \times \exp(\Sigma_t\ \Sigma_k\ \lambda_k f_k(y_{t-1},y_t,x,t))$$

- $Z(x)$: partition function
- $f_k$: feature functions
- $\lambda_k$: weights

## Belief Propagation

**Message Update**:

$$m_{i \to j}(x_j) = \Sigma_{xi}\ \psi_i(x_i)\psi_{ij}(x_i,x_j)\ \Pi_k{\in}N(i)\backslash j\ m_{k \to i}(x_i)$$

---

# More Probability Distributions

## Poisson Distribution

$$P(X = k) = (\lambda^k e^{-\lambda})\ /\ k!$$

- Mean = Variance = $\lambda$

## Exponential Distribution

$$f(x) = \lambda e^{-\lambda x}\ \text{for}\ x \geq 0$$

- Mean = $1/\lambda$
- Variance = $1/\lambda^2$

## Gamma Distribution

$$f(x) = (\beta^{\alpha}/\Gamma(\alpha)) \times x^{\alpha-1}e^{-\beta x}$$

## Beta Distribution

$$f(x) = (x^{\alpha-1}(1-x)^{\beta-1}) / B(\alpha,\beta)$$

where $B(\alpha,\beta) = \Gamma(\alpha)\Gamma(\beta)/\Gamma(\alpha+\beta)$

## Dirichlet Distribution

$$f(x_1,\ldots,x_k) = (1/B(\alpha)) \times \Pi_i \; x_i^{\alpha i-1}$$

## Student's t-Distribution

$$f(x) = \Gamma((v+1)/2) / (\sqrt{(v\pi)}\Gamma(v/2)) \times (1 + x^2/v)^{-(v+1)/2}$$

## Laplace Distribution

$$f(x) = (1/2b) \times \exp(-|x-\mu|/b)$$

---

# Anomaly Detection

## Gaussian Anomaly Detection

**Anomaly Score**:

$$f(x) = \exp(-\tfrac{1}{2}(x-\mu)^{\top}\Sigma^{-1}(x-\mu))$$

Anomaly if $f(x) < \varepsilon$

## Isolation Forest

**Path Length**:

$$h(x) = E[\text{path length to isolate } x]$$

**Anomaly Score**:

$$s(x,n) = 2^{-E(h(x))/c(n)}$$

## Local Outlier Factor (LOF)

```
LOF(x) = (Σᵧ∈N(x) lrd(y)) / (|N(x)| × lrd(x))
```

- lrd: local reachability density

## One-Class SVM

```
min (1/2)||w||² + (1/vn)Σᵢ ξᵢ - ρ
```

subject to: $w \cdot \varphi(x_i) \geq \rho - \xi_i$

## Mahalanobis Distance

```
D(x) = √((x-μ)ᵀS⁻¹(x-μ))
```

---

# Active Learning

## Uncertainty Sampling

**Least Confident**:

```
x* = argmax_x [1 - P(ŷ|x)]
```

**Margin Sampling**:

```
x* = argmin_x [P(ŷ₁|x) - P(ŷ₂|x)]
```

**Entropy-Based**:

```
x* = argmax_x [-Σᵢ P(yᵢ|x)log P(yᵢ|x)]
```

## Query by Committee

**Vote Entropy**:

```
x* = argmax_x [-Σᵢ (V(yᵢ)/C)log(V(yᵢ)/C)]
```

- $V(y_i)$: votes for class i
- C: committee size

## Expected Model Change

```
x* = argmax_x ||∇θL(θ;x,y)||
```

---

## Multi-Task Learning

### Hard Parameter Sharing Loss

```
L = Σt=1ᵀ αtLt(ft(x; θshared, θt), y⁽ᵗ⁾)
```

### Multi-Task Gaussian Process

```
f ~ GP(0, K ⊗ Kt)
```

- $K$: covariance between inputs
- $K_t$: covariance between tasks

### Task Uncertainty Weighting

```
L = Σt (1/2σt²)Lt + log σt
```

---

## Meta-Learning

### Model-Agnostic Meta-Learning (MAML)

**Inner Loop**:

```
θ'i = θ - α∇θLi(fθ)
```

**Outer Loop**:

```
θ = θ - β∇θ Σi Li(fθ'i)
```

### Prototypical Networks

**Prototype**:

```
ck = (1/|Sk|) × Σ(xi,yi)∈Sk fφ(xi)
```

**Prediction**:

```
P(y=k|x) = exp(-d(fφ(x),ck)) / Σk' exp(-d(fφ(x),ck'))
```

## Reptile Algorithm

$$\theta = \theta + \epsilon(\tilde{\theta} - \theta)$$

where $\tilde{\theta}$ is obtained after k steps of SGD

---

# Federated Learning

## FedAvg Algorithm

**Local Update**:

$$w_k^{t+1} = w_k^t - \eta\nabla F_k(w_k^t)$$

**Global Aggregation**:

$$w^{t+1} = \Sigma_{k=1}^K (n_k/n)w_k^{t+1}$$

## Differential Privacy in FL

$$\tilde{w} = w + N(0, \sigma^2 S^2 I)$$

- S: sensitivity
- σ: noise scale

---

# Fairness Metrics

## Demographic Parity

$$P(\hat{Y}=1|A=0) = P(\hat{Y}=1|A=1)$$

## Equalized Odds

$$P(\hat{Y}=1|A=0,Y=y) = P(\hat{Y}=1|A=1,Y=y) \text{ for } y\in\{0,1\}$$

## Equal Opportunity

$$P(\hat{Y}=1|A=0,Y=1) = P(\hat{Y}=1|A=1,Y=1)$$

## Disparate Impact

$$DI = P(\hat{Y}=1|A=0) / P(\hat{Y}=1|A=1)$$

Fair if DI ≥ 0.8

### Individual Fairness

$$d(f(x_1), f(x_2)) \leq Ld(x_1, x_2)$$

---

# Causal Inference

## Average Treatment Effect (ATE)

$$\tau = E[Y(1) - Y(0)]$$

## Propensity Score

$$e(x) = P(T=1|X=x)$$

## Inverse Probability Weighting

$$\hat{\tau}_{IPW} = (1/n)\Sigma_i [T_iY_i/e(X_i) - (1-T_i)Y_i/(1-e(X_i))]$$

## Doubly Robust Estimator

$$\hat{\tau}_{DR} = (1/n)\Sigma_i [\mu_1(X_i) - \mu_0(X_i) + T_i(Y_i-\mu_1(X_i))/e(X_i) - (1-T_i)(Y_i-\mu_0(X_i))/(1-e(X_i))]$$

## Instrumental Variables

$$\hat{\beta}_{IV} = Cov(Y,Z) / Cov(X,Z)$$

---

# Additional Important Concepts

## Rademacher Complexity

$$\hat{R}_n(F) = E_\sigma[\sup_{f\in F} (1/n)\Sigma_i \sigma_i f(x_i)]$$

- $\sigma_i$: Rademacher random variables

## VC Dimension

For hypothesis class H:

- VCD(H) = largest set size that can be shattered

## PAC Learning Bound

```
P(|R(h) - R̂(h)| > ε) ≤ 2exp(-2nε²)
```

## Margin Theory (SVM)

### Generalization Bound:

```
R(f) ≤ R̂_γ(f) + O(√(d/γ²n))
```

## Spectral Clustering

### Graph Laplacian:

```
L = D - W
```

### Normalized Laplacian:

```
L_norm = I - D⁻¹/²WD⁻¹/²
```

## Gaussian Processes

### Prior:

```
f(x) ~ GP(m(x), k(x,x'))
```

### Posterior:

```
f*|X,y,X* ~ N(μ*, Σ*)
```

where:

```
μ* = K*ᵀ(K+σ²I)⁻¹y
Σ* = K** - K*ᵀ(K+σ²I)⁻¹K*
```

## Variational Autoencoders (VAE)

### ELBO:

```
L = Eq(z|x)[log p(x|z)] - KL(q(z|x)||p(z))
```

## Generative Adversarial Networks (GAN)

**Objective**:

```
min_G max_D V(D,G) = Ex~pdata[log D(x)] + Ez~pz[log(1-D(G(z)))]
```

## Wasserstein Distance

```
W(P,Q) = inf_{γ∈Π(P,Q)} E(x,y)~γ[||x-y||]
```

---

# Modern Deep Learning Methods

## Diffusion Models

**Forward Process** (Adding Noise):

$$q(x_t|x_{t-1}) = N(x_t; \sqrt{(1-\beta_t)}x_{t-1}, \beta_t I)$$
$$q(x_t|x_0) = N(x_t; \sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)I)$$

where $\bar{\alpha}_t = \Pi_{s=1}^{t}\alpha_s$ and $\alpha_t = 1 - \beta_t$

**Reverse Process** (Denoising):

$$p\theta(x_{t-1}|x_t) = N(x_{t-1}; \mu\theta(x_t,t), \Sigma\theta(x_t,t))$$

**Training Objective** (Simplified):

$$L = E_t, x_0, \varepsilon[||\varepsilon - \varepsilon\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{(1-\bar{\alpha}_t)}\varepsilon, t)||^2]$$

**Score Matching**:

$$J(\theta) = \tfrac{1}{2} Ex,t[||s\theta(x,t) - \nabla x \log p_t(x)||^2]$$

**DDPM Sampling**:

$$x_{t-1} = 1/\sqrt{\alpha_t}(x_t - (1-\alpha_t)/\sqrt{(1-\bar{\alpha}_t)}\varepsilon\theta(x_t,t)) + \sigma_t z$$

## Normalizing Flows

**Change of Variables**:

$$\log p(x) = \log p(z) - \log|\det(\partial f/\partial x)|$$

**Flow Transformation**:

```
x = f(z) where z ~ p(z)
```

**Planar Flow**:

```
f(z) = z + uh(wᵀz + b)
```

$f(z) = z + uh(w^{T}z + b)$

**Real NVP Coupling Layer**:

```
y₁:d = x₁:d
yd+1:D = xd+1:D ⊙ exp(s(x₁:d)) + t(x₁:d)
```

$y_{1:d} = x_{1:d}$

$y_{d+1:D} = x_{d+1:D} \odot \exp(s(x_{1:d})) + t(x_{1:d})$

**Glow Objective**:

```
log p(x) = log p(z) + Σᵢ log|det(Jᵢ)|
```

$\log p(x) = \log p(z) + \sum_i \log|\det(J_i)|$

# Contrastive Learning

**InfoNCE Loss**:

$L = -\log\left[\exp(\text{sim}(z_i,z_j)/\tau) / \sum_{k=1}^{2N} \mathbb{1}[k\neq i]\exp(\text{sim}(z_i,z_k)/\tau)\right]$

- $\text{sim}(u,v) = u^{T}v/(\|u\|\cdot\|v\|)$
- $\tau$: temperature parameter

**SimCLR Objective**:

$l(i,j) = -\log\left[\exp(\text{sim}(z_i,z_j)/\tau) / \sum_{k=1}^{2N} \mathbb{1}[k\neq i]\exp(\text{sim}(z_i,z_k)/\tau)\right]$

**MoCo (Momentum Contrast)**:

$\theta_k \leftarrow m\theta_k + (1-m)\theta_q$

- m: momentum coefficient (e.g., 0.999)

**CLIP Loss**:

$L = -\frac{1}{2}\sum_i N\left[\log(\exp(\text{sim}(I_i,T_i)/\tau)/\sum_j \exp(\text{sim}(I_i,T_j)/\tau)) + \log(\exp(\text{sim}(I_i,T_i)/\tau)/\sum_j \exp(\text{sim}(I_j,T_i)/\tau))\right]$

## Vision Transformers (ViT)

**Patch Embedding**:

```
z₀ = [xclass; x¹pE; x²pE; ...; x^NpE] + Epos
```

- $x^p$: flattened patch
- E: embedding projection
- Epos: position embeddings

**Patch Size Calculation**:

```
N = HW/P² patches
```

- H,W: image height, width
- P: patch size

## Neural ODEs

**Continuous Dynamics**:

```
dh(t)/dt = f(h(t), t, θ)
```

**Forward Pass**:

```
h(t₁) = h(t₀) + ∫t₀^t¹ f(h(t), t, θ)dt
```

**Adjoint Sensitivity Method**:

```
da(t)/dt = -a(t)ᵀ ∂f/∂h
dL/dθ = -∫t₁^t₀ a(t)ᵀ ∂f/∂θ dt
```

---

# Advanced Reinforcement Learning

## Multi-Armed Bandits

### Upper Confidence Bound (UCB):

```
Aₜ = argmaxₐ[Q̂ₜ(a) + c√(ln t/Nₜ(a))]
```

- $N_t(a)$: times action a selected
- c: exploration constant

**Thompson Sampling**:

```
θₐ ~ Beta(αₐ, βₐ)
Aₜ = argmaxₐ θₐ
```

**ε-Greedy**:

```
Aₜ = {
    argmaxₐ Q̂ₜ(a)  with prob 1-ε
    random action   with prob ε
}
```

**Exp3 (Exponential-weight for Exploration)**:

```
pₜ(a) = (1-γ)wₜ(a)/Wₜ + γ/K
```

# Soft Actor-Critic (SAC)

## Soft Value Function:

```
V(s) = Eₐ~π[Q(s,a) - α log π(a|s)]
```

## Soft Q-Function:

```
Q(s,a) = r + γEₛ'[V(s')]
```

## Policy Objective:

```
J(π) = Eₛ~D,ε~N[α log π(f(ε;s)|s) - Q(s,f(ε;s))]
```

## Entropy Temperature Update:

```
α ← argminₐ Eₐ~π[-α log π(a|s) - αH̄]
```

# Trust Region Policy Optimization (TRPO)

## Surrogate Objective:

```
L(θ) = Eₛ,ₐ[πθ(a|s)/πθₒₗₐ(a|s) × A^π(s,a)]
```

## KL Constraint:

$$E_{s \sim \rho\pi}[KL(\pi\theta_{old}(\cdot|s)||\pi\theta(\cdot|s))] \leq \delta$$

**Natural Policy Gradient**:

$$\theta_{k+1} = \theta_k + \sqrt{(2\delta/g^{\mathsf{T}}Fg)} \times F^{-1}g$$

- F: Fisher information matrix
- g: policy gradient

## Proximal Policy Optimization (PPO)

**Clipped Objective**:

$$L(\theta) = E[\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1-\varepsilon, 1+\varepsilon)A_t)]$$

where $r_t(\theta) = \pi\theta(a_t|s_t)/\pi\theta_{old}(a_t|s_t)$

---

# Specialized Machine Learning Areas

## Survival Analysis

**Hazard Function**:

$$h(t) = \lim(\Delta t \to 0)\ P(t \leq T < t+\Delta t\ |\ T \geq t)/\Delta t$$

**Survival Function**:

$$S(t) = P(T > t) = \exp(-\int_0^t h(u)du)$$

**Cox Proportional Hazards**:

$$h(t|x) = h_0(t)\exp(\beta^{\mathsf{T}}x)$$

**Partial Likelihood**:

$$L(\beta) = \Pi_{i=1}^n\ [\exp(\beta^{\mathsf{T}}x_i) / \Sigma_{j \in R(t_i)}\ \exp(\beta^{\mathsf{T}}x_j)]^{\wedge}\delta_i$$

**Kaplan-Meier Estimator**:

$$\hat{S}(t) = \Pi_{i:t_i \leq t}\ [(n_i - d_i)/n_i]$$

## Learning to Rank

**RankNet Loss**:

```
L = -Σᵢⱼ P̄ᵢⱼlog Pᵢⱼ + (1-P̄ᵢⱼ)log(1-Pᵢⱼ)
```

where $P_{ij} = 1/(1 + \exp(-(s_i - s_j)))$

**LambdaRank Gradient**:

```
λᵢⱼ = -σ/(1 + exp(σ(sᵢ - sⱼ))) × |ΔNDCGᵢⱼ|
```

**ListNet Loss**:

```
L = -Σᵢ P(πᵢ|y)log P(πᵢ|s)
```

**NDCG (Normalized Discounted Cumulative Gain)**:

```
NDCG@k = DCG@k / IDCG@k
DCG@k = Σᵢ₌₁ᵏ (2^relᵢ - 1)/log₂(i + 1)
```

## Calibration Methods

**Platt Scaling**:

```
P(y=1|f(x)) = 1/(1 + exp(Af(x) + B))
```

**Temperature Scaling**:

```
q̂ᵢ = softmax(zᵢ/T)
```

**Isotonic Regression**: Fit monotonic function m: $\mathbb{R} \to$ [0,1]

```
min Σᵢ wᵢ(yᵢ - m(fᵢ))² s.t. m monotonic
```

**Expected Calibration Error (ECE)**:

```
ECE = Σₘ₌₁ᴹ |Bₘ|/n |acc(Bₘ) - conf(Bₘ)|
```

## Domain Adaptation

**Domain Adversarial Loss**:

```
L = Lᵧ(θf, θy) - λLd(θf, θd)
```

- $L_y$: task loss
- Ld: domain classifier loss

**Maximum Mean Discrepancy (MMD)**:

```
MMD²(X,Y) = ||𝔼[φ(x)] - 𝔼[φ(y)]||²ℋ
```

**CORAL Loss**:

```
L_CORAL = (1/4d²)||Cs - Ct||²F
```

- Cs, Ct: source/target covariance matrices

**Wasserstein Distance** (for DA):

```
W(Ps,Pt) = inf_{γ∈Π} 𝔼(xs,xt)~γ[||xs - xt||]
```

## Optimal Transport

### Kantorovich Problem:

```
W(μ,ν) = min_{π∈Π(μ,ν)} ∫ c(x,y)dπ(x,y)
```

### Sinkhorn Algorithm:

```
uⁱ⁺¹ = a ⊘ (Kvⁱ)
vⁱ⁺¹ = b ⊘ (Kᵀuⁱ⁺¹)
```

where K = exp(-C/ε)

### Wasserstein Barycenter:

```
min_{μ} Σₖ λₖW₂²(μ, μₖ)
```

### Gromov-Wasserstein Distance:

```
GW(C₁,C₂) = min_{π} Σᵢⱼₖₗ |C₁(i,k) - C₂(j,l)|²πᵢⱼπₖₗ
```

---

# Advanced Theoretical Concepts

# Conformal Prediction

**Prediction Set**:

```
C(x) = {y : s(x,y) ≥ q̂}
```

where $\hat{q}$ is $(1-\alpha)$ quantile of scores

**Coverage Guarantee**:

```
P(Y ∈ C(X)) ≥ 1 - α
```

**Conformity Score**:

```
s(x,y) = -|y - f(x)|   (regression)
s(x,y) = f(x)[y]       (classification)
```

# Advanced Differential Privacy

**Gaussian Mechanism**:

```
M(x) = f(x) + N(0, σ²S²f)
```

where $Sf = \max_{x,x'} \|f(x) - f(x')\|_2$

**Exponential Mechanism**:

```
P(M(x) = r) ∝ exp(εu(x,r)/(2Δu))
```

**Rényi Differential Privacy**:

```
Dα(P||Q) = (1/(α-1))log 𝔼q[(p(X)/q(X))^α]
```

**Moments Accountant**:

```
αM(λ) = max_aux log 𝔼[exp(λ·priv_loss(aux))]
```

# Advanced Kernel Methods

**String Kernel**:

```
K(s,t) = Σᵤ Σᵢ:ᵤ=ₛ[ᵢ] Σⱼ:ᵤ=ₜ[ⱼ] λ^(|i|+|j|)
```

**Graph Kernel** (Random Walk):

```
K(G₁,G₂) = Σᵢⱼ [I - λW×]⁻¹ᵢⱼ
```

**Fisher Kernel**:

```
K(x,x') = ∇θlog p(x|θ)ᵀF⁻¹∇θlog p(x'|θ)
```

**Polynomial Kernel with Offset**:

```
K(x,x') = (αxᵀx' + c)^d
```

## Capsule Networks

### Squashing Function:

```
vⱼ = ||sⱼ||²/(1 + ||sⱼ||²) × sⱼ/||sⱼ||
```

### Routing Algorithm:

```
cᵢⱼ = exp(bᵢⱼ) / Σₖ exp(bᵢₖ)
sⱼ = Σᵢ cᵢⱼûⱼ|ᵢ
```

### Agreement Update:

```
bᵢⱼ ← bᵢⱼ + ûⱼ|ᵢ · vⱼ
```

---

# Self-Supervised Learning

## Masked Autoencoders (MAE)

### Reconstruction Loss:

```
L = (1/|M|) Σᵢ∈M ||xᵢ - x̂ᵢ||²
```

- M: set of masked patches

**Masking Strategy**: Random masking with ratio r (typically 75%)

## BERT Objectives

### Masked Language Model (MLM):

```
L_MLM = -𝔼[Σᵢ∈M log P(xᵢ|x₍₋M₎)]
```

**Next Sentence Prediction (NSP)**:

```
L_NSP = -𝔼[log P(y|[CLS], A, B)]
```

**BERT Total Loss**:

```
L = L_MLM + L_NSP
```

# GPT Objectives

**Causal Language Modeling**:

$$L = -\sum_t \log P(x_t|x_1,\ldots,x_{t-1})$$

# Momentum Contrast (MoCo)

## Queue Update:

```
enqueue(queue, kq)
dequeue(queue)
```

## Momentum Update:

$$\theta_k \leftarrow m\theta_k + (1-m)\theta q$$

## InfoNCE Loss with Queue:

$$L = -\log[\exp(q\cdot k_+/\tau) / (\exp(q\cdot k_+/\tau) + \sum_{k-} \exp(q\cdot k_-/\tau))]$$

# BYOL (Bootstrap Your Own Latent)

**Loss** (No Negatives):

$$L = ||\bar{q} - z'||^2_2$$

where:

- $\bar{q} = q\theta(z\theta)$ (prediction)
- $z' = z\xi$ (target)

**Target Network Update**:

```
ξ ← τξ + (1-τ)θ
```

## SwAV (Swapping Assignments between Views)

**Clustering Assignment**:

```
Q = Sinkhorn(Z/ε)
```

**Swapped Prediction Loss**:

```
L = -Σᵢ[qᵢ⁽¹⁾log pᵢ⁽²⁾ + qᵢ⁽²⁾log pᵢ⁽¹⁾]
```

---

# Additional Modern Architectures

## EfficientNet Scaling

**Compound Scaling**:

```
depth: d = α^φ
width: w = β^φ
resolution: r = γ^φ
```

where $\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$

## Squeeze-and-Excitation

**Squeeze**:

```
zc = (1/HW) ΣᵢΣⱼ uc(i,j)
```

**Excitation**:

```
s = σ(W₂δ(W₁z))
```

**Scale**:

```
x̃c = sc · uc
```

## Feature Pyramid Networks

**Top-down Pathway**:

```
Pᵢ = Conv(Pᵢ₊₁↑ + Cᵢ)
```

## Focal Loss

```
FL(pₜ) = -αₜ(1-pₜ)^γ log(pₜ)
```

where:

```
pₜ = {
    p      if y = 1
    1-p    otherwise
}
```

## IoU Loss

### Standard IoU:

```
IoU = |A ∩ B| / |A ∪ B|
```

### GIoU (Generalized):

```
GIoU = IoU - |C\(A∪B)| / |C|
```

where C is smallest enclosing box

### DIoU (Distance):

```
DIoU = IoU - ρ²(b,b^gt) / c²
```

### CIoU (Complete):

```
CIoU = IoU - (ρ²(b,b^gt)/c² + αv)
```

where v measures aspect ratio consistency