# 1 Problem Formulation

**Input space:**

$\mathcal{X} \subseteq \mathbb{R}^d$

**Output space:**

$\mathcal{Y} = \{0,1\}$

**Data distribution:**
i.i.d. samples from unknown distribution.

**Learning Objective:**

$$\min_w \mathbb{E}[\ell(Y, f_w(X))]$$

# 2 Model Specification

**Hypothesis function:**

$$P(Y=1|X=x) = \sigma(w^T x)$$

Where sigmoid:

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

**Parameter space:**

$$w \in \mathbb{R}^d$$

**Structural assumptions:**

- Linear decision boundary
- Log-odds are linear in features

# 3 Loss Function

**Binary Cross-Entropy:**

$$L(y,\hat{y}) = -[y\log(\hat{y}) + (1-y)\log(1-\hat{y})]$$

**Why this loss?**

- Derived from Bernoulli likelihood
- Proper scoring rule

**Convexity:**

- Convex in $w$
- No closed-form solution

---

# 4 Objective Function

**Empirical Risk:**

$$J(w) = -\frac{1}{n}\sum_{i=1}^{n} \left[y_i \log(\sigma(w^T x_i)) + (1-y_i)\log(1-\sigma(w^T x_i))\right]$$

**Regularized (L2):**

$$J(w) + \lambda ||w||^2$$

---

# 5 Optimization Method

**Iterative method:**

- Gradient Descent
- Newton's Method

**Gradient:**

$$\nabla J(w) = \frac{1}{n} X^T (\sigma(Xw) - y)$$

**Convergence:**

- Convex $\Rightarrow$ global optimum
- Slower than linear regression

**Complexity:**

- $O(nd)O(nd)O(nd)$ per iteration

---

# 6 Statistical Interpretation

**MLE connection:**
Assume:

$Y|X \sim Bernoulli(p)Y|X \sim Bernoulli(p)Y|X \sim Bernoulli(p)$

Then maximizing likelihood = minimizing cross-entropy.

**Noise model:** Bernoulli

**Probabilistic meaning:**
Outputs:

$P(Y=1|X)P(Y=1|X)P(Y=1|X)$

---

# 7 Regularization & Generalization

**Bias–Variance:**

- High C (low $\lambda$) $\rightarrow$ overfitting
- High $\lambda$ $\rightarrow$ underfitting

**Capacity control:**

- L1 (feature selection)
- L2 (weight shrinkage)

---

# 8 Theoretical Properties

- Convex loss
- Unique global minimum
- Consistent classifier (under assumptions)

---

# 9 Computational Complexity

Training:

- $O(ndk)O(ndk)O(ndk)$ (k iterations)

Inference:

- $O(d)O(d)O(d)$

Memory:

- $O(nd)O(nd)O(nd)$

---

# 10 Limitations

- Cannot model non-linear boundaries (without feature engineering)
- Sensitive to multicollinearity
- Requires proper scaling
- Fails if classes perfectly separable (without regularization)