

---

# Categorical Distribution

## 1. What is the Categorical Distribution?

- **Generalises** the Bernoulli distribution from 2 categories to  $K$  categories.
- Models a *single* trial where the outcome is exactly one of  $K$  mutually-exclusive classes.
- Two common notations:
  - *Index form*: random variable  $Y \in \{1, \dots, K\}$ .
  - *One-hot form* ( $\mathbf{X}$ ): vector with exactly one 1 and the rest 0. **We use this form throughout.**

$$\mathbf{X} = (X_1, \dots, X_K)^\top, \quad X_i \in \{0, 1\}, \quad \sum_i X_i = 1.$$

Let  $\mathbf{p} = (p_1, \dots, p_K)^\top$  with  $p_i > 0$  and  $\sum_i p_i = 1$ .

## 2. Probability Mass Function (PMF)

$$P(\mathbf{X} = \mathbf{x}) = \prod_{i=1}^K p_i^{x_i}, \quad (\text{one-hot } \mathbf{x}).$$

**Why it works – “switch” intuition.** Because  $\mathbf{x}$  is one-hot, exactly one  $x_j = 1$  and all others = 0:

$$\prod_i p_i^{x_i} = p_j^1 \prod_{i \neq j} p_i^0 = p_j,$$

so the product “switches on” the sole probability matching the 1 entry—identical in spirit to the Bernoulli PMF when  $K = 2$ .

## 3. Worked Example

Suppose  $K = 3$  categories (“red”, “blue”, “green”):

$$\mathbf{p} = \begin{pmatrix} 0.2 \\ 0.6 \\ 0.2 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad (\text{blue chosen}).$$

Then

$$P(\mathbf{X} = \mathbf{x}) = 0.2^0 0.6^1 0.2^0 = 0.6.$$

---

## 4. Expectation (Mean Vector)

$$\mathbb{E}[\mathbf{X}] = \sum_{\mathbf{x}} \mathbf{x} P(\mathbf{X} = \mathbf{x}) = \mathbf{p}.$$

*Interpretation:* over many trials, the “hot 1” spends a proportion  $p_i$  of the time in slot  $i$ .

## 5. Covariance Matrix

$$\Sigma = \text{Cov}(\mathbf{X}) = \mathbb{E}[\mathbf{X}\mathbf{X}^\top] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{X}]^\top = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top.$$

### Break-down.

- $\mathbb{E}[\mathbf{X}\mathbf{X}^\top]$ : For a one-hot  $\mathbf{X}$  this is a matrix with a single 1 on the diagonal position that was chosen. Taking expectation over all possibilities  $\Rightarrow \text{diag}(\mathbf{p})$ .
- Off-diagonal entries are *negative*. When one category turns “on”, all others must be 0, so they move in opposite directions.

$$\Sigma = \begin{pmatrix} p_1(1-p_1) & -p_1p_2 & \dots & -p_1p_K \\ -p_2p_1 & p_2(1-p_2) & \dots & -p_2p_K \\ \vdots & \vdots & \ddots & \vdots \\ -p_Kp_1 & -p_Kp_2 & \dots & p_K(1-p_K) \end{pmatrix}.$$

- Symmetric; every row and column sums to 0.
- Variance of each category appears on the diagonal; covariances (always  $\leq 0$ ) off the diagonal.

## 6. Special Case – Bernoulli Inside

If  $K = 2$  and  $\mathbf{p} = (p, 1 - p)$  then

$$P(\mathbf{X} = \mathbf{x}) = p^{x_1}(1-p)^{x_2},$$

which is exactly the Bernoulli PMF written in one-hot form.

## 7. Where is it used?

- **Classification** targets (one-hot labels).
- **Generative models** (e.g. predicting next token).
- **Reinforcement learning**: policy selects an action among  $K$  possibilities.

---

## 8. Summary

- Categorical( $\mathbf{p}$ ) models a single draw from  $K$  classes.
- PMF:  $\prod_i p_i^{x_i}$ .
- Mean vector:  $\mathbf{p}$ .
- Covariance:  $\text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top$ .
- Reduces to Bernoulli when  $K = 2$ .