

# Convolutional Neural Networks & Classification

## Data Science in Electron Microscopy

---

Philipp Pelz

2024

[https://github.com/ECLIPSE-Lab/SS24\\_DataScienceForEM](https://github.com/ECLIPSE-Lab/SS24_DataScienceForEM)

# Regression

- Used to answer *how much?* or *how many?* questions.
  - Predicting house prices.
  - Estimating baseball team wins.
  - Forecasting patient hospital stay duration.
- Important distinctions within regression:
  - House price: Non-negative, often relative to baseline price.
    - Regress on logarithm of price.
  - Hospital stay duration: Discrete, nonnegative.
    - Least mean squares may not be ideal.
    - Specialized subfield: *survival modeling*.

# Key Points

- Estimation involves more than minimizing squared errors.
- Supervised learning encompasses more than regression.

## Focus of this Section

- Classification problems.
  - Shift focus from *how much?* to *which category?* questions.

## Classification Examples

- Does this email belong in the spam folder or the inbox?
- Is this customer more likely to sign up or not to sign up for a subscription service?
- Does this image depict a donkey, a dog, a cat, or a rooster?
- Which movie is Aston most likely to watch next?
- Which section of the book are you going to read next?

# Understanding Classification

- Colloquially, *classification* refers to:
  1. Hard assignments of examples to categories (classes).
  2. Soft assignments, assessing the probability of each category.
- Distinction between hard and soft assignments often blurred:
  - Models making soft assignments used even for hard assignment tasks.

## Multi-Label Classification

- Some cases involve more than one true label:
  - Example: A news article covering entertainment, business, and space flight, but not medicine or sports.
- Known as **multi-label classification**(<https://en.wikipedia.org/wiki/>)

# Simple Image Classification Problem

## Problem Setup

- **Input:** grayscale image.
  - Each pixel value represented by a scalar.
  - Four features:  $x_1, x_2, x_3, x_4$ .
- **Categories:** "cat", "chicken", "dog".

## Label Representation

## Natural Impulse

- Use for respectively.
  - Efficient for storage:  $y \in \{\text{dog, cat, chicken}\}$
  - Suitable for ordinal data with natural ordering, e.g., age categories.

## Ordinal Regression

- Useful for naturally ordered categories.
- Resources:
  - Overview of ranking loss functions: :citet:Moon.Smola.Chang.ea.2010.
  - Bayesian approach for multi-modal responses: :citet:Beutel.Murray.Faloutsos.ea.2014.

# One-Hot Encoding

- Suitable for unordered categorical data.
- Vector with components for each category.
- Example:
  - “cat” →  $(1, 0, 0)$
  - “chicken” →  $(0, 1, 0)$
  - “dog” →  $(0, 0, 1)$

$$y \in \{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}.$$

# Linear Model

## Estimating Conditional Probabilities

- **Objective:** Estimate conditional probabilities for all classes.
- **Model Requirement:** Multiple outputs (one per class).

## Affine Functions for Classification

- Need as many affine functions as output categories.
- For symmetry, use redundant parametrization:
  - Last category derived from others.
- **Case Example:**
  - 4 features, 3 output categories.
  - 12 scalars for weights () and 3 scalars for biases ().

$$w_{ij} \qquad b_i$$

# Model Equations

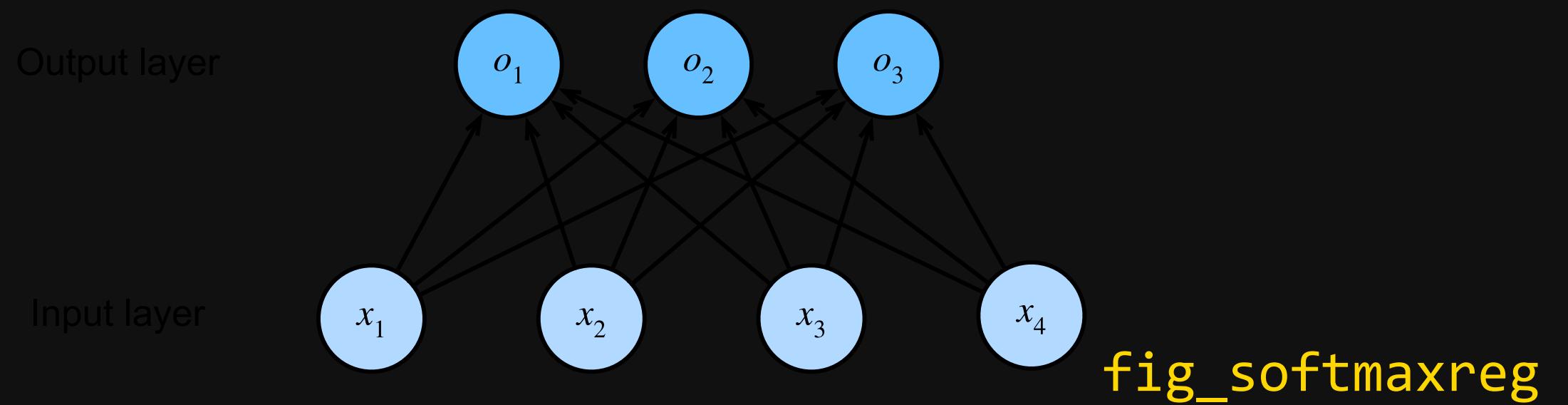
$$o_1 = x_1 w_{11} + x_2 w_{12} + x_3 w_{13} + x_4 w_{14} + b_1,$$

$$o_2 = x_1 w_{21} + x_2 w_{22} + x_3 w_{23} + x_4 w_{24} + b_2,$$

$$o_3 = x_1 w_{31} + x_2 w_{32} + x_3 w_{33} + x_4 w_{34} + b_3.$$

# Neural Network Diagram

- **Diagram:** See :numref:[fig\\_softmaxreg](#).
- **Structure:**
  - Single-layer neural network (similar to linear regression).
  - Each output () depends on all inputs ( $x_1, x_2, \dots$ ).
  - Output layer described as a *fully connected layer*.



# Concise Notation with Vectors and Matrices

- Use vectors and matrices for concise notation:

- 

- $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$  Suited for mathematics and code implementation.

- **Weights and Biases:**

- Weights gathered into a matrix.

$$3 \times 4$$

- Biases in a vector.

$$\mathbf{b} \in \mathbb{R}^3$$

# The Softmax

:label:subsec\_softmax\_operation

## Limitations of Vector-Valued Regression

- Directly minimizing the difference between  $\hat{y}$  and labels  $y$  has drawbacks:
  - No guarantee that outputs sum to 1 (expected for probabilities).
  - Outputs might be negative or exceed 1 even if they sum to 1.
- These issues complicate the estimation problem and make solutions brittle to outliers.

## Example Problem

- Positive linear dependency between number of bedrooms and likelihood of buying a house:
  - Probability might exceed 1 for a mansion.

# Solution: Squish the Outputs

## Probit Model

- Assume outputs  $\mathbf{o}$  are corrupted versions of :

■ , with .  $\mathbf{o} \sim \mathcal{N}(0, \sigma^2)$

■ Known as the **probit model**, introduced by :[citet:Fechner.1860](#).

■ Less effective and harder to optimize compared to softmax.

## Softmax Function

- Use an exponential function: .

■ Ensures nonnegative probabilities.  $P(y = i) \propto \exp o_i$

■ Monotonic relationship: Conditional class probability increases with .

- Normalize to ensure probabilities sum to 1:

■ Divide each exponential value by their sum.

# Softmax Function

:eqlabel: eq\_softmax\_y\_and\_o

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{o}) \quad \text{where} \quad \hat{y}_i = \frac{\exp(o_i)}{\sum_j \exp(o_j)}.$$

- The largest coordinate of  $\hat{\mathbf{y}}$  corresponds to the most likely class according to  $\hat{y}_i$ .
- Softmax preserves the ordering among its arguments:
  - To determine the highest probability class, compute:

$$\operatorname{argmax}_j \hat{y}_j = \operatorname{argmax}_j o_j.$$

# Vectorization

## Computational Efficiency with Vectorization

- **Minibatch:**
  - examples, each with inputs.
- **Output Categories:** categories
  - Weights:  $\mathbf{W} \in \mathbb{R}^{d \times q}$
  - Bias:  $\mathbf{b} \in \mathbb{R}^{1 \times q}$

# Vectorized Equations

eq\_minibatch\_softmax\_reg

$$\mathbf{O} = \mathbf{X}\mathbf{W} + \mathbf{b},$$
$$\hat{\mathbf{Y}} = \text{softmax}(\mathbf{O}).$$

- **Efficiency:** Converts the dominant operation to matrix-matrix product .  
 $\mathbf{X}\mathbf{W}$
- **Rowwise Softmax:**
  - Exponentiate entries in each row of .  
 $\mathbf{O}$
  - Normalize by the sum for each row.

# Numerical Stability

- **Care Needed:** Avoid overflow/underflow with large numbers.
- **Deep Learning Frameworks:** Handle these issues automatically.

# Loss Function

## Mapping Features to Probabilities

- Goal: Optimize the accuracy of mapping from features to probabilities .
- Method: Maximum likelihood estimation.
  - Same concept as used for probabilistic justification of mean squared error loss in :numref:[subsec\\_normal\\_distribution\\_and\\_squared\\_loss](#).

$$\mathbf{x} \qquad \hat{\mathbf{y}}$$

# Log-Likelihood

- **Softmax Output:** Vector interpreted as estimated conditional probabilities:
  - Example:  $\hat{\mathbf{y}}$
- **Dataset Representation:**  $\hat{y}_1 = P(y_1 = \text{cat} | \mathbf{x})$ 
  - Features:
  - Labels: (one-hot encoded).
- **Probability Comparison:**
  - Compare estimated probabilities with actual classes:

$$P(\mathbf{Y} | \mathbf{X}) = \prod_{i=1}^n P(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}).$$

- **Factorization:** Assumes each label is drawn independently

# Negative Log-Likelihood

- $\log P(\mathbf{Y} \mid \mathbf{X}) = \sum_{i=1}^n -\log P(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}) = \sum_{i=1}^n l(\mathbf{y}$
- **Loss Function:** For any pair of label  $\mathbf{y}$  and model prediction  $\hat{\mathbf{y}}$  over classes  $\hat{q}$

$$l(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{j=1}^q y_j \log \hat{y}_j.$$

:eqlabel: eq\_1\_cross\_entropy

# Cross-Entropy Loss

- **Definition:** Loss function in :eqref: eq\_1\_cross\_entropy is called *cross-entropy loss*.
- **One-Hot Vector:** is a one-hot vector of length  $q$ .
  - Sum over all coordinates  $j$  vanishes for all but one term.

## Properties of Cross-Entropy Loss

- **Bounded Below by 0:** Loss  $l(\mathbf{y}, \hat{\mathbf{y}})$  is bounded from below by 0 when  $\mathbf{y}$  is a probability vector.
  - No single entry larger than 1, negative logarithm cannot be lower than 0.
  - only if the actual label is predicted with certainty.
  - Predicting with certainty requires input to go to infinity (or others to negative infinity).
- **Infinite Loss:** Assigning a probability of 0 incurs infinite loss ().
  - High confidence errors result in infinite loss.  $-\log 0 = \infty$

# Softmax and Cross-Entropy Loss

:label:subsec\_softmax\_and\_derivatives

## Calculation of Softmax and Cross-Entropy Loss

- **Combining Equations:**

- Plugging :eqref:eq\_softmax\_y\_and\_o into the definition of the loss in :eqref:eq\_1\_cross\_entropy:

$$= \sum_{j=1}^q y_j \log \sum_{k=1}^q \exp(o_k) - \sum_{j=1}^q y_j o_j$$

$$= \log \sum_{k=1}^q \exp(o_k) - \sum_{j=1}^q y_j o_j.$$

# Derivative with Respect to Logit

$o_j$

$$\partial_{o_j} l(\mathbf{y}, \hat{\mathbf{y}}) = \frac{\exp(o_j)}{\sum_{k=1}^q \exp(o_k)} - y_j = \text{softmax}(\mathbf{o})_j - y_j$$

- **Interpretation:**

- Derivative is the difference between model probability (softmax) and actual label (one-hot vector).
- Similar to regression where gradient was the difference between observation and estimate .
- In exponential family models, gradients of log-likelihood are given by this term  $y_j - \hat{y}_j$ , making gradient computation easy.

# Distribution Over Outcomes

- **General Case:**
  - Observe entire distribution over outcomes.
  - Label  $\hat{y}$  as a probability vector (e.g.,  $(0.1, 0.2, 0.7)$ ) instead of binary vector.
- **Mathematical Consistency:**
  - The same loss function in :eqref:`eq\_1\_cross\_entropy` applies.
  - Represents expected value of the loss for a distribution over labels.

## Cross-Entropy Loss

# Information Theory Basics

:label:subsec\_info\_theory\_basics

Many deep learning papers use intuition and terms from information theory. To make sense of them, we need some common language. This is a survival guide. *Information theory* deals with the problem of encoding, decoding, transmitting, and manipulating information (also known as data).

## Entropy

- **Central Idea:** Quantify the amount of information contained in data, placing a limit on data compression.
- **Definition:** For a distribution  $P$ , its *entropy* is:

$$P$$

eq\_softmax\_reg\_entropy

$$H[P] = \sum_j -P(j) \log P(j).$$

- **Fundamental Theorem:** To encode data drawn randomly from distribution  $P$ , we need at least “nats” to encode it :cite:Shannon.1948.  $H[P]$
- **Nat:** Equivalent of a bit but using base  $e$  instead of base 2.
- $1 \text{ nat} \approx 1.44 \text{ bits.}$

## Surprisal

- **Compression and Prediction:** Easy-to-predict data is easy to compress.
- Example: A data stream where every token is the same is easy to predict and compress.
- **Surprise:** Greater when an event has a lower assigned probability.
- Claude Shannon quantified *surprisal* at observing an event with probability  $P(j)$  as  $\log \frac{1}{P(j)}$ .
- **Expected Surprisal:** Entropy defined in :eqref:softmax\_reg\_entropy is the  $\sum_j P(j) \log \frac{1}{P(j)}$  when probabilities match the data-generating process.

# Cross-Entropy Revisited

## Definition and Interpretation

- **Entropy:** The level of surprise experienced by someone who knows the true probability.
- **Cross-Entropy:** The expected surprisal of an observer with subjective probabilities upon seeing data generated according to probabilities .

$P$

- Denoted as .
- Formula: .  $H(P, Q)$
- Lowest cross-entropy achieved when  $H(P, Q) \stackrel{\text{def}}{=} \sum_{j} P(j) \log Q(j)$

$$H(P, Q) = H(P)$$

## Cross-Entropy in Classification

- Two perspectives on cross-entropy classification objective:
  1. Maximizing the likelihood of observed data.
  2. Minimizing surprisal (and thus the number of bits) required to communicate the labels.

# Summary and Discussion

## Key Points

- Introduced the first nontrivial loss function for optimizing over *discrete* output spaces.
- Used a probabilistic approach, treating discrete categories as instances of draws from a probability distribution.
- Encountered the softmax activation function, transforming neural network outputs into valid discrete probability distributions.
- Derivative of cross-entropy loss combined with softmax is similar to the derivative of squared error (difference between expected behavior and prediction).
- Briefly touched on connections to statistical physics and information theory.

# Exercises

1. We can explore the connection between exponential families and the softmax in some more depth.

1. Compute the second derivative of the cross-entropy loss for the softmax.

2. Compute the variance of the distribution given by  $\text{softmax}(\mathbf{o})$  and show that it matches the second derivative computed above.

2. Assume that we have three classes which occur with equal probability, i.e., the probability vector is .

1. What is the problem if we try to design a binary code for it?

2. Can you design a better code? Hint: what happens if we try to encode two independent observations? What if we encode observations jointly?

3. When encoding signals transmitted over a physical wire, engineers do not always use binary codes.

For instance, PAM-3 uses three signal levels as opposed to two levels. How many ternary units do you need to transmit an integer in the range  $\{0, \dots, 7\}$ ? Why might this be a better idea in terms of electronics?



1. The Bradley-Terry model uses a logistic model to capture preferences. For a user to choose between apples and oranges one assumes scores  $a$  and  $b$ . Our requirements are that larger scores should lead to a higher likelihood in choosing the associated item and that the item with the largest score is the most likely one to be chosen :cite:Bradley.Terry.1952.

1. Prove that the softmax satisfies this requirement.
2. What happens if you want to allow for a default option of choosing neither apples nor oranges?

Hint: now the user has 3 choices.

2. Softmax derives its name from the following mapping: .

$$\text{RealSoftMax}(a, b) = \log(\exp(a) + \exp(b))$$

1. Prove that .
2. How small can you make the difference between both functions? Hint: without loss of generality you can set  $a = 0$  and  $b = 1$ .  
 $\text{RealSoftMax}(a, b) > \max(a, b)$
3. Prove that this holds for  $a > b$ , provided that .  
 $b = a - \lambda^{-1} \text{RealSoftMax}(a, b)$
4. Show that for we have .  
 $\lambda^{-1} \text{RealSoftMax}(\lambda a, \lambda b) \rightarrow \max(a, b)$
5. What does the soft-min look like?  
 $\lambda^{-1} \text{RealSoftMin}(a, b) \rightarrow \min(a, b)$
6. Extend this to more than two numbers.

1. The function is sometimes also referred to as the log-partition function.

1. Prove that the function  $\hat{g}(\mathbf{x}) \stackrel{\text{def}}{=} \log \sum_i \exp x_i$  is convex. Hint: to do so, use the fact that the first derivative amounts to the probabilities from the softmax function and show that the second derivative is the variance.

2. Show that is translation invariant, i.e., .

3. What happens if some of the coordinates  $a$  are very large? What happens if they're all very small?

4. Show that if we choose we end up with a numerically stable implementation.

2. Assume that we have some probability distribution . Suppose we pick another distribution with for .

1. Which choice of corresponds to doubling the temperature? Which choice corresponds to halving it?  
 $b = \max_i x_i$

2. What happens if we let the temperature converge to ?

3. What happens if we let the temperature converge to ?

$\infty$

# The Image Classification Dataset

## MNIST Dataset

- **Overview:** Widely used for image classification of handwritten digits.
  - **Reference:** MNIST dataset :cite:LeCun.Bottou.Bengio.ea.1998.
  - **Content:** 60,000 images of  $28 \times 28$  pixels resolution (plus a test dataset of 10,000 images).
- **Historical Context:**
  - Released in the 1990s, it was a formidable challenge for machine learning algorithms.
  - Example of state-of-the-art equipment at the time: Sun SPARCStation 5 with 64MB RAM and 5 MFLOPs (AT&T Bell Laboratories, 1995).
- **Impact:**
  - High accuracy on digit recognition was crucial for automating letter sorting for the USPS in the 1990s.
  - Algorithms achieving error rates below 1%:
    - Deep networks like LeNet-5 :cite:LeCun.Jackel.Bottou.ea.1995.
    - Support vector machines with invariances :cite:Scholkopf.Burges.Vapnik.1996.
    - Tangent distance classifiers :cite:Simard.LeCun.Denker.ea.1998.

# MNIST as a Benchmark

- **Historical Significance:** Served as the reference for comparing machine learning algorithms for over a decade.
- **Current Relevance:**
  - Simple models achieve classification accuracy over 95%.
  - Unsuitable for distinguishing between stronger and weaker models due to high levels of accuracy.
  - Skewed algorithmic development towards methods that excel with clean datasets.

## Transition to Fashion-MNIST

- **Limitations of MNIST:** More of a sanity check than a benchmark today.
- **ImageNet:**
  - Poses a more relevant challenge :cite:Deng.Dong.Socher.ea.2009.
  - Too large for interactive examples and illustrations in this context.
- **Fashion-MNIST:**
  - Released in 2017 :cite:Xiao.Rasul.Vollgraf.2017.
  - Contains images of 10 categories of clothing at  $28 \times 28$  pixels resolution.
  - Used as a substitute for MNIST in upcoming sections.

```
1 %matplotlib inline
2 import time
3 from d2l import torch as d2l
4 import torch
5 import torchvision
6 from torchvision import transforms
7
8 d2l.use_svg_display()
```

## Loading the Dataset

Since it is such a frequently used dataset, all major frameworks provide preprocessed versions of it. We can [download and read the Fashion-MNIST dataset into memory using built-in framework utilities.]

```
1 class FashionMNIST(d2l.DataModule): #@save
2     """The Fashion-MNIST dataset."""
3     def __init__(self, batch_size=64, resize=(28, 28)):
4         super().__init__()
```

```
5     self.save_hyperparameters()
6     trans = transforms.Compose([transforms.Resize(resize),
7                                 transforms.ToTensor()])
8     self.train = torchvision.datasets.FashionMNIST(
9         root=self.root, train=True, transform=trans, download=True)
10    self.val = torchvision.datasets.FashionMNIST(
11        root=self.root, train=False, transform=trans, download=True)
```

## Fashion-MNIST Dataset

- **Content:** Images from 10 categories.
  - Each category represented by:
    - 6,000 images in the training dataset.
    - 1,000 images in the test dataset.
- **Dataset Sizes:**
  - **Training Set:** 60,000 images.
  - **Test Set:** 10,000 images.
- **Purpose of Test Dataset:** Used exclusively for evaluating model performance, not for training.

```
1 data = FashionMNIST(resize=(32, 32))
2 len(data.train), len(data.val)
```

## Image Characteristics

- **Grayscale Images:** Upscaled to pixels resolution.
  - Similar to original MNIST dataset ( $32 \times 32$  binary black and white images).
- **Modern Image Data:**
  - Typically has 3 channels (red, green, blue).
  - Hyperspectral images can have over 100 channels (e.g., HyMap sensor has 126 channels).
- **Storage Convention:**
  - Images stored as a tensor.
  - $c \times h \times w$  = number of color channels.
  - $c$  = height.
  - $h$  = width.

```
1 data.train[0][0].shape
```

The categories of Fashion-MNIST have human-understandable names. The following convenience method converts between numeric labels and their names.

```
1 @d2l.add_to_class(FashionMNIST) #@save
2 def text_labels(self, indices):
3     """Return text labels."""
4     labels = ['t-shirt', 'trouser', 'pullover', 'dress', 'coat',
5               'sandal', 'shirt', 'sneaker', 'bag', 'ankle boot']
6     return [labels[int(i)] for i in indices]
```

## Reading a Minibatch

- **Data Iterator:** Utilizes built-in data iterator for reading from training and test sets.
- **Minibatch Size:** At each iteration, a data iterator reads a minibatch of data with size `batch_size`.
- **Shuffling:** Randomly shuffle the examples for the training data iterator to ensure better training performance.

```
1 @d2l.add_to_class(FashionMNIST) #@save
2 def get_dataloader(self, train):
3     data = self.train if train else self.val
4     return torch.utils.data.DataLoader(data, self.batch_size, shuffle=train,
5                                         num_workers=self.num_workers)
```

To see how this works, let's load a minibatch of images by invoking the `train_dataloader` method. It contains 64 images.

```
1 x, y = next(iter(data.train_dataloader()))
2 print(x.shape, x.dtype, y.shape, y.dtype)
```

## Performance of Data Loading

- **Time to Read Images:** Not extremely fast, but sufficient for our needs.
- **Context:** Processing images with a deep network takes significantly longer than loading them.
- **Conclusion:** The data loader's speed ensures that training a network will not be IO constrained.

```
1 tic = time.time()
2 for x, y in data.train_dataloader():
3     continue
4 f'{time.time() - tic:.2f} sec'
```

## Visualization

- **Frequent Use:** We'll use the Fashion-MNIST dataset frequently.
- **Convenience Function:** `show_images` function to visualize images and labels.
  - Implementation details in the appendix.

```
1 def show_images(imgs, num_rows, num_cols, titles=None, scale=1.5):    #@save
2     """Plot a list of images."""
3     raise NotImplementedError
4
5     @d2l.add_to_class(FashionMNIST)    #@save
6     def visualize(self, batch, nrows=1, ncols=8, labels=[]):
7         X, y = batch
8         if not labels:
9             labels = self.text_labels(y)
10        if tab.selected('mxnet', 'pytorch'):
11            d2l.show_images(X.squeeze(1), nrows, ncols, titles=labels)
12        if tab.selected('tensorflow'):
13            d2l.show_images(tf.squeeze(X), nrows, ncols, titles=labels)
14        if tab.selected('jax'):
15            d2l.show_images(jnp.squeeze(X), nrows, ncols, titles=labels)
16
17    batch = next(iter(data.val_dataloader()))
18    data.visualize(batch)
```

We are now ready to work with the Fashion-MNIST dataset in the sections that follow.

## Summary

# The Base Classification Model

:label:sec\_classification

- **Observation:** Implementations from scratch and concise implementations using frameworks are similar for both regression and classification.
- **Purpose:** Since many models in this book deal with classification, we aim to add functionalities to support this setting specifically.
- **Base Class:** This section provides a base class for classification models to simplify future code.

```
1 from d2l import torch as d2l  
2 import torch
```

# The Classifier Class

- **Definition:** `Classifier` class.
- **Validation Step:**
  - Reports both the loss value and classification accuracy on a validation batch.
  - Update every `num_val_batches` batches.
  - Generates averaged loss and accuracy on the whole validation data.
  - Average numbers may not be exact if the last batch contains fewer examples, but this is ignored for simplicity.

```
1 class Classifier(d2l.Module): #@save
2     """The base class of classification models."""
3     def validation_step(self, batch):
4         Y_hat = self(*batch[:-1])
5         self.plot('loss', self.loss(Y_hat, batch[-1]), train=False)
6         self.plot('acc', self.accuracy(Y_hat, batch[-1]), train=False)
```

By default we use a stochastic gradient descent optimizer, operating on minibatches, just as we did in the context of linear regression.

```
1 @d2l.add_to_class(d2l.Module) #@save
2 def configure_optimizers(self):
3     return torch.optim.SGD(self.parameters(), lr=self.lr)
```

# Accuracy

- **Predicted Probability Distribution:**  $y_{\text{hat}}$ 
  - Choose the class with the highest predicted probability for hard predictions.
  - Example: Gmail categorizing emails into “Primary”, “Social”, “Updates”, “Forums”, or “Spam”.
- **Correct Predictions:** When predictions match the label class  $y$ .
  - **Classification Accuracy:** Fraction of correct predictions.
  - Difficult to optimize directly as it is not differentiable, but crucial for performance measurement and benchmarks.
  - Almost always reported when training classifiers.

- **Computing Accuracy:**

- Assume  $y_{\text{hat}}$  is a matrix with prediction scores for each class in the second dimension.
- Use `argmax` to get the predicted class index for the largest entry in each row.
- Compare predicted class with ground-truth  $y$  elementwise.
- Convert  $y_{\text{hat}}$ 's data type to match  $y$  to ensure correct comparisons.
- Result: Tensor with 0 (false) and 1 (true) entries.
- Sum the tensor to get the number of correct predictions.

```
1 @d2l.add_to_class(Classifier) #@save
2 def accuracy(self, Y_hat, Y, averaged=True):
3     """Compute the number of correct predictions."""
4     Y_hat = d2l.reshape(Y_hat, (-1, Y_hat.shape[-1]))
5     preds = d2l.astype(d2l.argmax(Y_hat, axis=1), Y.dtype)
6     compare = d2l.astype(preds == d2l.reshape(Y, -1), d2l.float32)
7     return d2l.reduce_mean(compare) if averaged else compare
```

# Summary

- **Common Problem:** Classification is common enough to warrant its own convenience functions.
- **Central Importance:** *Accuracy* of the classifier.
  - Often the primary measure of performance.
- **Training Objectives:**
  - Train classifiers to optimize various objectives for statistical and computational reasons.
  - Regardless of the loss function minimized during training, assessing accuracy empirically is crucial.
- **Convenience Method:** Useful to have methods for assessing classifier accuracy.

# Exercises

- Denote by  $L_v$  the validation loss, and let  $\bar{L}_v^q$  be its quick and dirty estimate computed by the loss function averaging in this section. Lastly, denote  $\bar{L}_v^b$  by the loss on the last minibatch. Express  $\bar{L}_v^q$  in terms of  $L_v$ ,  $\bar{L}_v^b$ , and the sample and minibatch sizes.
- Show that the quick and dirty estimate  $\bar{L}_v^q$  is unbiased. That is, show that  $E[\bar{L}_v^q] = E[L_v]$ . Why would you still want to use  $\bar{L}_v^q$  instead?
- Given a multiclass classification loss, denoting by  $l$  the penalty of estimating  $y'$  when we see  $y$  and given a probability  $p(y | x)$ , formulate the rule for an optimal selection of  $y'$ . Hint: express the expected loss, using  $p(y | x)$  and  $l$ .

# Softmax Regression Implementation from Scratch

## Implementing Softmax Regression

- **Fundamental Concept:** Softmax regression is a foundational model.
- **Importance:** Essential to understand and implement it yourself.
- **Approach:**
  - Define softmax-specific aspects of the model.
  - Reuse components from the linear regression section, including the training loop.

```
1 from d2l import torch as d2l  
2 import torch
```

# The Softmax

- **Key Concept:** Mapping from scalars to probabilities.
- **Sum Operator:** Recall the operation of the sum operator along specific dimensions in a tensor, as discussed in :numref:`subsec\_lin-alg-reduction` and :numref:`subsec\_lin-alg-non-reduction`.
- **Example:**
  - Given a matrix **X**, sum over all elements by default or over elements in the same axis.
  - The **axis** variable allows for computing row and column sums.

```
1 X = d2l.tensor([[1.0, 2.0, 3.0], [4.0, 5.0, 6.0]])
2 d2l.reduce_sum(X, 0, keepdims=True), d2l.reduce_sum(X, 1, keepdims=True)
```

- **Computing the Softmax:** Requires three steps:
  1. Exponentiation of each term.
  2. Sum over each row to compute the normalization constant for each example.
  3. Division of each row by its normalization constant, ensuring the result sums to 1.

$$\text{softmax}(\mathbf{X})_{ij} = \frac{\exp(\mathbf{X}_{ij})}{\sum_k \exp(\mathbf{X}_{ik})}.$$

- **Partition Function:**

- The (logarithm of the) denominator is called the (log) *partition function*.
- Introduced in **statistical physics** to sum over all possible states in a thermodynamic ensemble.

- **Implementation:** The implementation is straightforward.

```
1 def softmax(X):  
2     X_exp = d2l.exp(X)  
3     partition = d2l.reduce_sum(X_exp, 1, keepdims=True)  
4     return X_exp / partition ## The broadcasting mechanism is applied here
```

- **Input Transformation:** For any input  $X$ , [each element is turned into a non-negative number, and each row sums up to 1], as required for a probability.
- **Caution:**
  - The code above is *not* robust against very large or very small arguments.
  - While sufficient to illustrate the concept, do *not* use this code verbatim for serious purposes.
- **Best Practice:**
  - Deep learning frameworks have built-in protections.
  - Use the built-in softmax functions in deep learning frameworks going forward.

```
1 x = d2l.rand((2, 5))
2 x_prob = softmax(X)
3 x_prob, d2l.reduce_sum(x_prob, 1)
```

# The Model

- **Softmax Regression Model:** We now have everything needed to implement the softmax regression model.
- **Instance Representation:**
  - Each instance represented by a fixed-length vector.
  - Raw data consists of pixel images.
  - [Flatten each image $^{28 \times 28}$ , treating them as vectors of length 784.]
- **Future Improvements:**
  - Later chapters will introduce convolutional neural networks to exploit spatial structure more effectively.
- **Outputs:**
  - Number of outputs equals the number of classes.
  - (Dataset has 10 classes, so the network has an output dimension of 10.)

- **Weights and Biases:**

- Weights: matrix.
- Biases:  $784 \times 10$  dimensional row vector.
- Initialize  $1 \times 10$  weights  $\mathbf{W}$  with Gaussian noise.
- Initialize biases as zeros.

```
1
2 class SoftmaxRegressionScratch(d2l.Classifier):
3     def __init__(self, num_inputs, num_outputs, lr, sigma=0.01):
4         super().__init__()
5         self.save_hyperparameters()
6         self.W = torch.normal(0, sigma, size=(num_inputs, num_outputs),
7                               requires_grad=True)
8         self.b = torch.zeros(num_outputs, requires_grad=True)
9
10    def parameters(self):
11        return [self.W, self.b]
```

# Model Implementation

- **Mapping Input to Output:**

- The code below defines the network's mapping from input to output.
- Each pixel image in the batch is flattened into a vector using `reshape` before passing through the model  $28 \times 28$

```
1 class SoftmaxRegression(nn.Module):  
2     def __init__(self, num_inputs, num_outputs):  
3         super(SoftmaxRegression, self).__init__()  
4         self.linear = nn.Linear(num_inputs, num_outputs)  
5  
6     def forward(self, X):  
7         X = X.reshape((-1, 784)) # Flatten each image  
8         return self.linear(X)  
  
1 @d2l.add_to_class(SoftmaxRegressionScratch)  
2 def forward(self, X):  
3     X = d2l.reshape(X, (-1, self.W.shape[0]))  
4     return softmax(d2l.matmul(X, self.W) + self.b)
```

# The Cross-Entropy Loss

- **Common Loss Function:**

- The cross-entropy loss function is perhaps the most common in deep learning.
- Widely used in classification problems, more so than regression problems.

- **Definition:**

- Cross-entropy takes the negative log-likelihood of the predicted probability assigned to the true label.
- Avoid Python for-loops for efficiency; use indexing instead.
- One-hot encoding in `np` allows selecting matching terms in  $\hat{y}$ .

- **Example:**

$$\mathbf{y} \quad \hat{\mathbf{y}}$$

- Create sample data `y_hat` with 2 examples of predicted probabilities over 3 classes and their corresponding labels `y`.
- Correct labels: `y = [0, 2]` (i.e., the first and third class).
- Use `y` as indices of probabilities in `y_hat` to pick out terms efficiently.

```
1 y = d2l.tensor([0, 2])
2 y_hat = d2l.tensor([[0.1, 0.3, 0.6], [0.3, 0.2, 0.5]])
3 y_hat[[0, 1], y]
```

50

Now we can (**implement the cross-entropy loss function**) by averaging over the logarithms of the selected probabilities.



# Training

- **Reusing fit Method:** Use the `fit` method defined in :numref:`sec\_linear\_scratch` to [train the model with 10 epochs.]
- **Adjustable Hyperparameters:**
  - **Number of epochs** (`max_epochs`)
  - **Minibatch size** (`batch_size`)
  - **Learning rate** (`lr`)
  - These hyperparameters are not learned during the primary training loop but influence model performance.
- **Hyperparameter Tuning:**
  - Choose values based on the *validation* split of the data.
  - Evaluate the final model on the *test* split.
- **Evaluation:**
  - Treat the test data of Fashion-MNIST as the validation set.
  - Report validation loss and validation accuracy on this split.

```

1 data = d2l.FashionMNIST(batch_size=256)
2 model = SoftmaxRegressionScratch(num_inputs=784, num_outputs=10, lr=0.1)
3 trainer = d2l.Trainer(max_epochs=10)
4 trainer.fit(model, data)

```

## Prediction

Now that training is complete, our model is ready to [classify some images.]

```

1 X, y = next(iter(data.val_dataloader()))
2 if tab.selected('pytorch', 'mxnet', 'tensorflow'):
3     preds = d2l.argmax(model(X), axis=1)
4 if tab.selected('jax'):
5     preds = d2l.argmax(model.apply({'params': trainer.state.params}, X), axis=1)
6 preds.shape

```

## Error Analysis

- **Focus:** More interested in images labeled *incorrectly*.
- **Visualization:**
  - Compare actual labels (first line of text output) with model predictions (second line of text output).

```

1 wrong = d2l.astype(preds, y.dtype) != y
2 X, y, preds = X[wrong], y[wrong], preds[wrong]
3 labels = [a+'\n'+b for a, b in zip(
4     data.text_labels(y), data.text_labels(preds))]
5 data.visualize([X, y], labels=labels)

```

# Summary

- **Experience Gained:**
  - Solving linear regression and classification problems.
  - Reached the state of the art of 1960-1970s statistical modeling.
- **Next Steps:**
  - Learn how to leverage deep learning frameworks.
  - Implement models more efficiently using these frameworks.

# Exercises

1. In this section, we directly implemented the softmax function based on the mathematical definition of the softmax operation. As discussed in :numref:sec\_softmax this can cause numerical instabilities.

1. Test whether softmax still works correctly if an input has a value of ?
  2. Test whether softmax still works correctly if the largest of all inputs is  $^{100}$  smaller than ?
  3. Implement a fix by looking at the value relative to the largest entry in the argument. $^{-100}$
2. Implement a cross\_entropy function that follows the definition of the cross-entropy loss function .

1. Try it out in the code example above.
2. Why do you think it runs more slowly?
3. Should you use it? In which cases would it make sense?
4. What do you need to be careful of? Hint: consider the domain of the logarithm.
3. Is it always a good idea to return the most likely label? For example, would you do this for medical diagnosis? How would you try to address this?
4. Assume that we want to use softmax regression to predict the next word based on some features. What are some problems that might arise from a large vocabulary?
5. Experiment with the hyperparameters of the code above. In particular:
  1. Plot how the validation loss changes as you change the learning rate.
  2. Do the validation and training loss change as you change the minibatch size? How large or small do you need to go before you see an effect?



# Generalization and Overfitting

- **Focus:** Tackling multiclass classification with linear neural networks, softmax functions, and cross-entropy loss.
  - Interpreted model outputs as probabilistic predictions.
  - Derived the cross-entropy loss function (negative log likelihood).
  - Fitted the model to the training set.
- **Goal:** Learn *general patterns* assessed empirically on previously unseen data (the test set).
  - High accuracy on the training set is meaningless if it only reflects memorization.
- **Problem with Memorization:**
  - Perfect accuracy on the training set can be achieved by memorizing the dataset.
  - Memorization doesn't help classify new examples.
  - Without further guidance, new examples might require random guessing.

- **Key Questions:**
  1. How many test examples are needed to precisely estimate the accuracy of classifiers on the underlying population?
  2. What happens if we keep evaluating models on the same test set repeatedly?
  3. Why should fitting our linear models to the training set fare better than naive memorization?

# Overfitting and Generalization

- **Previous Coverage:**
  - Basics of overfitting and generalization discussed in :numref:`sec\_generalization\_basics` in the context of linear regression.
- **Current Chapter:**
  - Delve deeper into foundational ideas of statistical learning theory.
  - Explore guarantees of generalization *a priori*.
- **Generalization Guarantees:**
  - For many models, we can guarantee generalization within an upper bound  $\epsilon$  on the generalization gap.
  - Determine the required number of samples  $n$  to ensure the empirical error lies within  $\epsilon$  of the true error, for any data generating distribution.

- **Practical Utility:**
  - These guarantees are profound but limited in practical use for deep learning practitioners.
  - Ensuring generalization of deep neural networks *a priori* would require an absurd number of examples (potentially trillions).
  - In practice, deep neural networks generalize well with far fewer examples (thousands).
- **Practical Approach:**
  - Forgo *a priori* guarantees.
  - Employ methods based on past performance on similar problems.
  - Certify generalization *post hoc* through empirical evaluations.
- **Future Discussion:**
  - Revisit generalization in :numref:chap\_perceptrons.
  - Provide an introduction to the scientific literature explaining why deep neural networks generalize in practice.

# The Test Set

- **Role of Test Sets:**
  - Used as the gold standard method for assessing generalization error.
- **Focus:**
  - Discuss properties of generalization error estimates using a fixed classifier .
  - Assume a *fresh* dataset of examples not used to train the classifier.  $f$
- **Empirical Error:**
  - Fraction of instances where the prediction disagrees with the true label .
  - Given by:  $\epsilon_{\mathcal{D}}(f) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(f(\mathbf{x}^{(i)}) \neq y^{(i)}).$

$$\epsilon_{\mathcal{D}}(f) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(f(\mathbf{x}^{(i)}) \neq y^{(i)}).$$

- **Population Error:**

- Expected fraction of examples in the underlying population (distribution  $P(X, Y)$ ) where the classifier disagrees with the true label.
- Given by:

$$\epsilon(f) = E_{(\mathbf{x}, y) \sim P} \mathbf{1}(f(\mathbf{x}) \neq y) = \int \int \mathbf{1}(f(\mathbf{x}) \neq y) p(:$$

# Estimating Population Error

- **Population Error:**
  - is the actual quantity of interest but cannot be directly observed.
  - $\epsilon(f)$  similar to estimating the average height in a large population without measuring every individual.
- **Empirical Error as Estimator:**
  - Test set  $\mathcal{D}$  is statistically representative of the population.
  - serves as a statistical estimator of  $\epsilon(f)$

- **Mean Estimation:**

- Estimating population error is a classic problem of mean estimation (see :numref:sec\_prob).
- Central limit theorem: With random samples drawn from any distribution with mean  $\mu$  and standard deviation  $\sigma$ , the sample average approximates a normal distribution centered at  $\hat{\mu}$  with standard deviation  $\sigma/\sqrt{n}$ .

- **Convergence Rate:**  $\sigma/\sqrt{n}$

- As the number of examples grows, test error approaches true error at a rate of  $\mathcal{O}(1/\sqrt{n})$ .
- To estimate test error twice as precisely, collect four times as large a test set.
- To reduce test error by a factor of one hundred, collect ten thousand times as large a test set.

- **Bernoulli Random Variable:**

- Random variable can take values 0 and 1 (Bernoulli random variable).
- Parameter indicating probability of value 1 is the true error rate .  
 $\epsilon(f)$
- Variance of a Bernoulli distribution is .  
 $\sigma^2 = \epsilon(f)(1 - \epsilon(f))$
- Variance is highest when the true error rate is close to 0.5 and lower when close to 0 or 1.

- **Asymptotic Standard Deviation:**

- Asymptotic standard deviation of over test samples cannot be greater than .

$$\epsilon_{\mathcal{D}}(f) \sqrt{0.25/n}$$

# Estimating Population Error (Continued)

- **Finite Sample Considerations:**

- Asymptotic rate characterizes behavior as test set size approaches infinity.
- For  $\epsilon_D(f)$  to approximate  $\epsilon(f)$  within  $\pm 0.01$  (one standard deviation), collect roughly 2500 samples.
- For  $\epsilon_D(f)$  to approximate  $\epsilon(f)$  within  $\pm 0.01$  ( $\pm 0.01$  confidence interval); collect 10000 samples.

- **Popular Benchmarks:**  $\epsilon_D(f) \in \epsilon(f) \pm 0.01$

- Test set sizes for many popular benchmarks in machine learning are around 10000 samples.
- Improvements of  $\pm 0.01$  in error rates are significant, especially when error rates are close to  $0$ .

- **Finite Sample Bounds:**

- Asymptotic analysis provides general insights but for finite samples, use Hoeffding's inequality:

$$P(\epsilon_{\mathcal{D}}(f) - \epsilon(f) \geq t) < \exp(-2nt^2).$$

- **Example Calculation:**

- To conclude with 95% confidence that  $\epsilon_{\mathcal{D}}(f)$  is within  $t$  of  $\epsilon(f)$ , roughly 15000 examples are needed (compared to 10000 from asymptotic analysis).  $\epsilon_{\mathcal{D}}(f) = 0.01(f)$

- **General Trend:**

- Finite sample guarantees are more conservative but not drastically different from asymptotic estimates.
  - Asymptotic analysis remains useful for providing ballpark figures.



# Test Set Reuse

- **Empirical ML Research:**
  - Models developed and validated based on test set performance.
  - Evaluate test error for any fixed classifier to understand population error .
- **Training Initial Model**  $f_1$   $f$   $\epsilon(f)$ 
  - Train model with an appropriate number of test examples.
  - Use a validation set for preliminary analysis, hyperparameter tuning, and model selection.
  - Evaluate on the test set, report unbiased population error estimate with confidence interval.

- **Developing New Models:**

- Develop a new model and tune it on the validation set.
- Prepare to evaluate but realize the test set has already been used.

- **Issues with Test Set Reuse:**

- Original test set was designed for a single classifier .
- Evaluating multiple classifiers on the same test set introduces the problem of false discovery.
- With multiple classifiers, the probability of at least one misleading result increases.
- Example: With 20 classifiers, ensuring no misleading score is challenging.

- **Multiple Hypothesis Testing:**

- Evaluating multiple hypotheses on the same dataset can lead to misleading conclusions.
- Persistent problem in scientific research despite extensive statistical literature.



# Concerns with Test Set Reuse

- **Distrusting Subsequent Evaluations:**
  - Initial analysis assumes the classifier was chosen without contact with the test set.
  - Testing multiple functions and using the test set performance of  $f_1$  to choose  $f_2$  introduces bias.
- **Adaptive Overfitting:**
  - Occurs when information from the test set leaks to the modeler.
  - Once leaked, the test set can no longer be a true test set.
  - Topic of interest in learning theory and statistics :cite:dwork2015preserving.

- **Practical Guidance:**
  - Create and use real test sets.
  - Consult test sets infrequently.
  - Account for multiple hypothesis testing when reporting confidence intervals.
  - Increase vigilance when stakes are high and dataset size is small.
- **Benchmark Challenges:**
  - Maintain several test sets.
  - Demote old test sets to validation sets after each round.

# Statistical Learning Theory

- **Importance of Test Sets:**
  - *Test sets are all that we really have* for evaluating generalization.
  - Seldom possess a true test set; often, someone else has already used it.
  - Even with a true test set, evaluating subsequent models can be frustrating and unreliable.
- **Limitations of Test Sets:**
  - Only provide *post hoc* evidence of generalization.
  - Do not provide *a priori* reasons to expect generalization.
- **Appeal of Statistical Learning Theory:**
  - Aims to elucidate principles explaining why/when models trained on empirical data generalize to unseen data.
  - Primary aim: bound the generalization gap, relating model properties to dataset size.

- **Empirical vs. True Error:**

- Learning theorists seek to bound the difference between the *empirical error* on training set  $\mathcal{S}$  and the *true error* on the population.  
 $\epsilon_{\mathcal{S}}(f_{\mathcal{S}})$
- Different from evaluation of a fixed classifier on a separate test set.

- **Challenges with Classifier Collections:**

- Easy to estimate error for a single fixed classifier.
- Difficult to estimate error accurately when considering collections of classifiers.
- High probability of at least one classifier in the collection having a misleadingly low error.

- **Function Class Considerations:**

- When choosing a classifier from a set of functions  $\mathcal{F}$ , there's a risk of selecting one with a grossly underestimated population error.
- Linear models with continuously valued parameters typically involve choosing from an infinite class of functions  $\mathcal{F}$ .  
 $|\mathcal{F}| = \infty$

# Statistical Learning Theory (Continued)

- **Uniform Convergence:**
  - Ambitious solution: develop tools for proving uniform convergence.
  - Goal: Empirical error rate for every classifier in the class  $\mathcal{F}$  converges to its true error rate with high probability.
  - Theoretical principle: With probability at least  $1 - \delta$ , no classifier's error rate  $\epsilon(f)$  will be misestimated by more than some small amount .
- **Model Class Flexibility:**
  - Not all model classes  $\mathcal{F}$  can achieve uniform convergence.
  - Example: Class of memorizers (always achieve empirical error 0 but fail on the population).
  - Tradeoff: Flexible (high variance) models fit training data well but risk overfitting. Rigid (high bias) models generalize well but risk underfitting.

- **Central Question:**
  - Quantify where a model sits on the bias-variance spectrum.
  - Provide guarantees on model performance.
- **Vapnik-Chervonenkis (VC) Dimension:**
  - Series of seminal papers by Vapnik and Chervonenkis extended theory on convergence of relative frequencies.
  - VC dimension measures the complexity (flexibility) of a model class.
  - Key result: Bound on the difference between empirical error and population error as a function of VC dimension and number of samples:

$$P(R[p, f] - R_{\text{emp}}[\mathbf{X}, \mathbf{Y}, f] < \alpha) \geq 1 - \delta \text{ for } \alpha \geq$$

# Statistical Learning Theory (Continued)

- **Key Parameters:**

- $\epsilon$ : Probability that the bound is violated.
- $\delta > 0$ : Upper bound on the generalization gap.
- $\alpha$ : Dataset size.
- $n$ : Constant depending on the scale of the loss.

- **Using the Bound:**

- Plug in desired values of  $\epsilon$  and  $\alpha$  to determine the required number of samples.

- **VC Dimension:**

- Measures the largest number of data points for which any arbitrary (binary) labeling can be assigned and find a model  $f$  in the class that agrees with that labeling.
- Example: Linear models on  $d$ -dimensional inputs have VC dimension  $d+1$ .
  - A line can assign any possible labeling to three points in two dimensions, but not to four.

- **Practical Considerations:**

- The theory can be overly pessimistic for complex models.
- Obtaining the guarantee typically requires more examples than actually needed to achieve the desired error rate.
- Fixing the model class and  $\delta$ , the error rate decays with the usual rate.

- **Generalization Gap:**

$$\mathcal{O}(1/\sqrt{n})$$

- While we may not do better in terms of  $n$ , varying the model class can present a pessimistic picture of the generalization gap.

# Summary

- **Model Evaluation:**
  - Use a test set of previously unseen data for evaluation.
  - Test set evaluations provide an unbiased estimate of true error.
  - Converge at the  $O(1/\sqrt{n})$  rate as the test set grows.
  - Provide approximate confidence intervals based on asymptotic distributions or conservative finite sample guarantees.
- **Challenges with Test Sets:**
  - Often not true test sets (used repeatedly by multiple researchers).
  - Multiple evaluations on the same test set make controlling for false discovery difficult.
  - Significance of the problem depends on holdout set size and usage.
  - Good practice: curate real test sets (or multiple) and use them conservatively.

- **Statistical Learning Theory:**
  - Developed methods for guaranteeing uniform convergence over a model class.
  - Uniform convergence ensures that every model's empirical error converges to its true error simultaneously.
  - Allows for choosing the model that minimizes training error with confidence in its generalization.
- **VC Dimension:**
  - Introduced by Vladimir Vapnik and Alexey Chervonenkis.
  - Measures the complexity of a model class.
  - Provides uniform convergence results for all models in a VC class.
  - Training errors for all models in the class grow closer to their true errors at rates.  
 $\mathcal{O}(1/\sqrt{n})$

- **Alternative Complexity Measures:**
  - Numerous measures proposed for generalization guarantees (see :citet:boucheron2005theory for detailed discussion).
  - Useful in statistical theory but less effective for explaining deep neural networks' generalization.
- **Deep Neural Networks:**
  - Often have millions of parameters and can assign random labels to large datasets.
  - Generalize well on practical problems, often better when larger and deeper, despite larger VC dimensions.
  - Next chapter will revisit generalization in the context of deep learning.

# Exercises

1. If we wish to estimate the error of a fixed model to within  $\epsilon$  with probability greater than 99.9%, how many samples do we need?  $f = 0.0001$
2. Suppose that somebody else possesses a labeled test set  $\mathcal{D}$  and only makes available the unlabeled inputs (features). Now suppose that you can only access the test set labels by running a model (no restrictions placed on the model class) on each of the unlabeled inputs and receiving the corresponding error  $\epsilon$ . How many models would you need to evaluate before you leak the entire test set and thus could appear to have error  $\epsilon$ , regardless of your true error?
3. What is the VC dimension of the class of  $0^{\text{th}}$ -order polynomials?
4. What is the VC dimension of axis-aligned  $5^{\text{th}}$  rectangles on two-dimensional data?

# Environment and Distribution Shift

:label:sec\_environment-and-distribution-shift

- **Context:**

- Previous sections focused on fitting models to datasets.
- Did not consider where data originates or the ultimate use of model outputs.

- **Common Pitfall:**

- Rushing to develop models without considering data origins and usage.
- Leads to failed deployments when data distribution shifts.

- **Distribution Shift Issues:**

- Models perform well on test sets but fail in deployment.
- Deployment can perturb data distribution.
- Example: Loan repayment model using footwear as a feature might cause behavior change without improving credit-worthiness.

- **Critical Considerations:**

- Decision-making based on model outputs can alter the environment.
- Awareness of such impacts is crucial for responsible use of machine learning.

- **Mitigating Distribution Shift:**

- Expose common concerns and stimulate critical thinking.
- Detect situations early and mitigate damage.
- Use machine learning responsibly.

- **Solutions:**

- Simple: Ask for the “right” data.
- Technically difficult: Implement reinforcement learning systems.
- Ethical: Address philosophical questions about algorithm application.

# Types of Distribution Shift

- **Focus:**
  - Passive prediction setting.
  - Ways data distributions might shift.
  - Salvaging model performance.
- **Classic Setup:**
  - Training data sampled from distribution  $p_S(\mathbf{x}, y)$
  - Test data drawn from a different distribution  $p_T(\mathbf{x}, y)$
  - Learning a robust classifier is impossible without assumptions on how  $p_S$  and  $p_T$  relate.

- **Example:**

- Binary classification (dogs vs. cats).
- If  $p_S(y|T(x)) = 1 - p_T(y|x)$ , labels are flipped without changing the input distribution.
- Impossible to distinguish from a setting where the distribution did not change.

- **Restricted Assumptions:**

- Under certain assumptions, algorithms can detect shifts.
- Algorithms can sometimes adapt on the fly, improving classifier accuracy.

# Covariate Shift

- **Definition:**

- Covariate shift occurs when the distribution of inputs changes over time, but the labeling function ( $P(y|x)$ ) does not change.
- The problem arises due to a shift in the distribution of the covariates (features).

- **Causality:**

- Covariate shift is a natural assumption when  $x$  causes  $y$ .

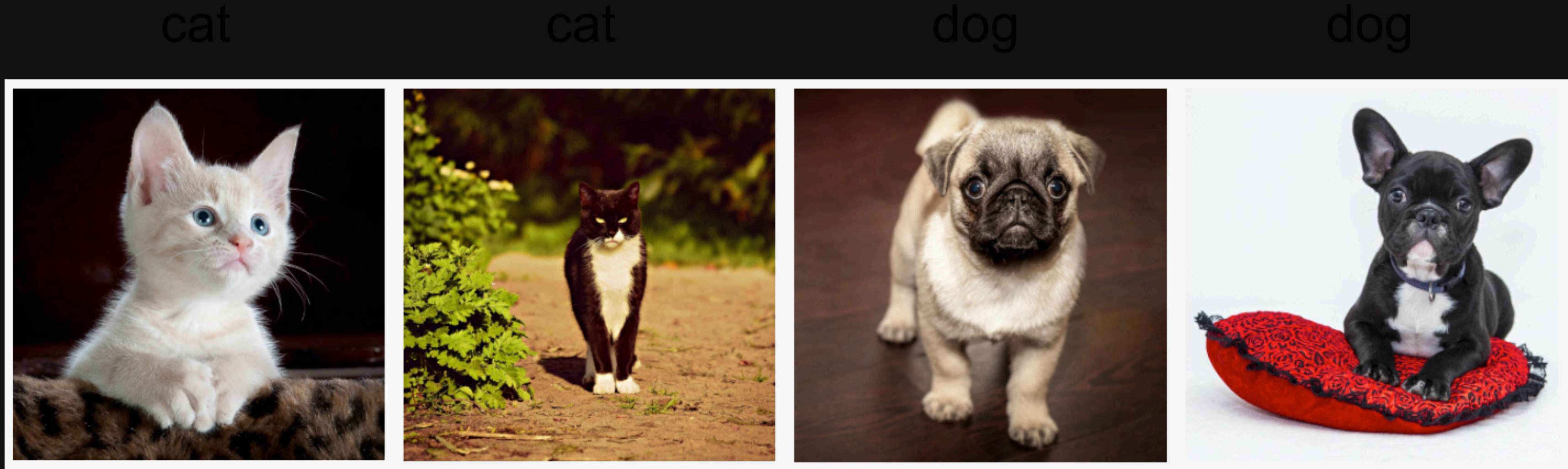
$x$        $y$



- **Example:** Distinguishing cats and dogs.

- **Training Data:**

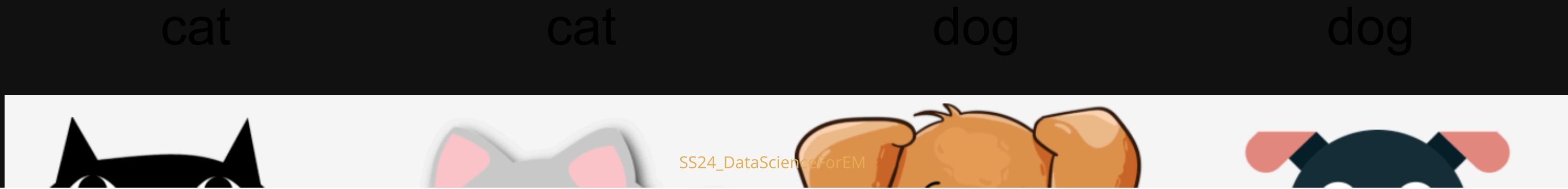
- Consists of photos (see :numref:**fig\_cat-dog-train**).



:label:**fig\_cat-dog-train**

- **Test Data:**

- Consists of cartoons (see :numref:**fig\_cat-dog-test**).



# Label Shift

- **Definition:**

- Label shift occurs when the label marginal  $P(y)$  can change, but the class-conditional distribution  $P(\mathbf{x} | y)$  remains fixed across domains.
- Assumption to make when  $y$  causes  $\mathbf{x}$ .

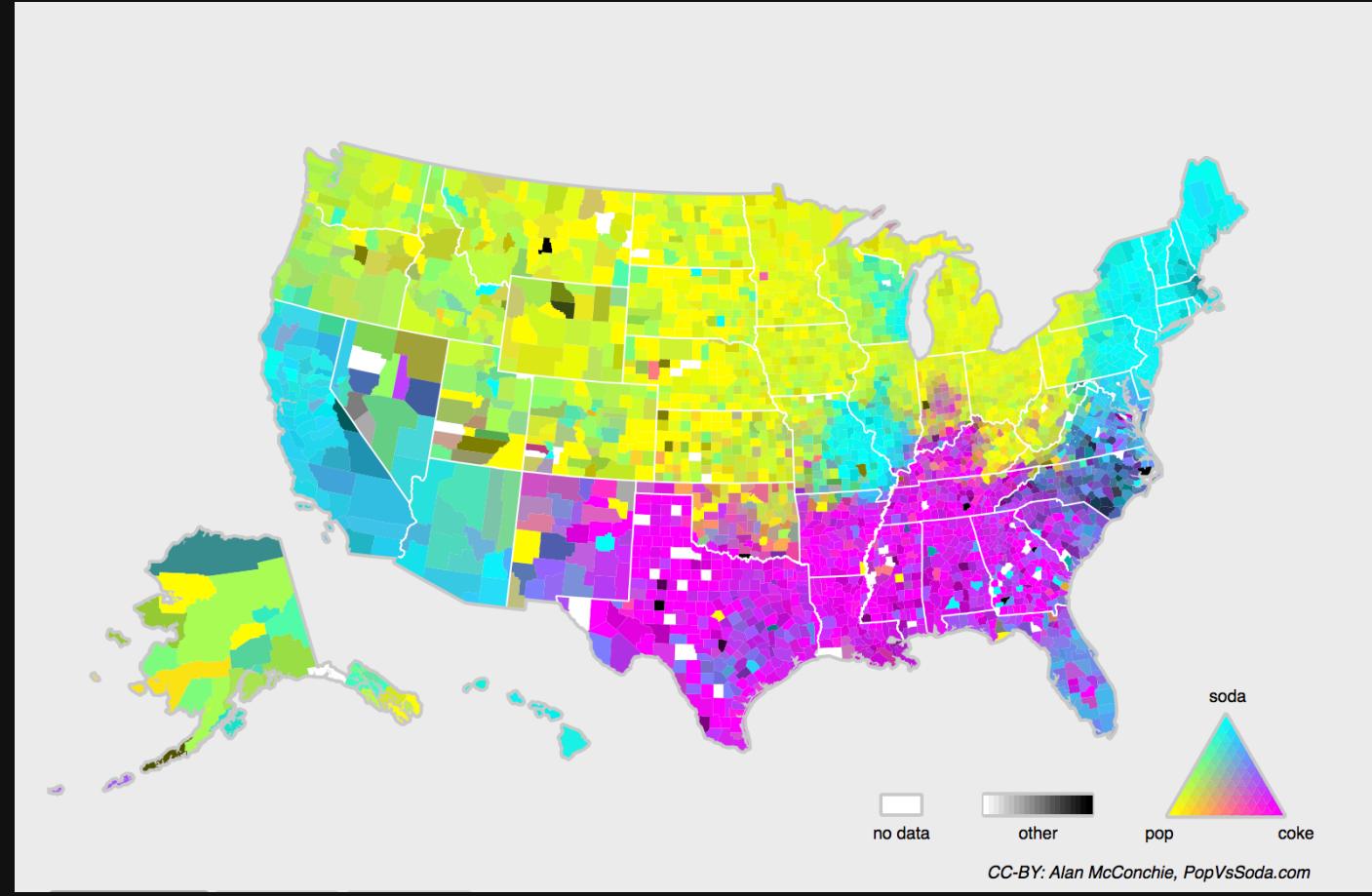
- **Example:**

- Predicting diagnoses given symptoms.
- Relative prevalence of diagnoses changes over time, but diseases cause symptoms.

- **Simultaneous Occurrence:**
  - Label shift and covariate shift can hold simultaneously in degenerate cases.
  - When the label is deterministic, covariate shift assumption is satisfied even when  $y$  causes  $\mathbf{x}$ .
- **Methodological Advantage:**
  - Methods based on label shift assumption often involve manipulating low-dimensional objects (labels).
  - This is advantageous in deep learning where inputs are high-dimensional.

# Concept Shift

- **Definition:**
  - Concept shift occurs when the definitions of labels change over time or across different contexts.
  - Example: Categories like diagnostic criteria for mental illness, fashion trends, and job titles.
- **Example:**
  - Geographic variation in names for *soft drinks* in the United States (see :numref:[fig\\_popvssoda](#)).



fig\_popvssoda

- **Challenges:**
  - can vary depending on location (e.g., machine translation systems).
  - $P(y|x)$  shift can be subtle and hard to detect.
  - Shift often occurs gradually, either temporally or geographically.

## Examples of Distribution Shift

Before delving into formalism and algorithms, let's discuss some concrete situations where covariate or concept shift might not be obvious.

# Nonstationary Distributions

- **Definition:**
  - Distribution changes slowly over time (nonstationary).
  - Model is not updated frequently enough.
- **Examples:**
  - Computational advertising model not updated with new devices (e.g., iPad launch).
  - Spam filter that fails to adapt to new spam tactics.
  - Product recommendation system recommending seasonal items out of season (e.g., Santa hats after Christmas).

# Correction of Distribution Shift

- **Overview:**
  - Training and test distributions often differ.
  - Sometimes models work despite distribution shifts; other times, principled strategies are needed.

## Empirical Risk and Risk

:label:subsec\_empirical-risk-and-risk

- **Model Training:**
  - Iterate over features and labels of training data .
  - Update model parameters after every minibatch.  $\{(x_1, y_1), \dots, (x_n, y_n)\}$
- **Empirical Risk Minimization:**
  - Minimize the loss on the training data:

:eqlabel: eq\_empirical-risk-min

$$\underset{f}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n l(f(\mathbf{x}_i), y_i),$$

- : Loss function measuring “how bad” the prediction is given the label .
- *Empirical risk*: Average loss over the training data to approximate the *task*.
- *Risk*: Expectation of the loss over the entire population of data drawn from true distribution .

$$p(\mathbf{x}, y)$$

:eqlabel: eq\_true-risk

$$E_{p(\mathbf{x}, y)}[l(f(\mathbf{x}), y)] = \int \int l(f(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} dy.$$

- **Practical Strategy:**

- In practice, we cannot obtain the entire population of data.
- *Empirical risk minimization* is used to approximate minimizing the risk.

# Covariate Shift Correction

- **Goal:** Estimate using labeled data drawn from a *source distribution* instead of the *target distribution*.
  - **Dependency Assumption:** Conditional distribution does not change
  - **Risk Correction:**
- $P(y_i \mid \mathbf{x}_i)$        $(\mathbf{x}_i, y_i)$        $q(\mathbf{x})$        $p(\mathbf{x})$   
 $p(y \mid \mathbf{x}) = q(y \mid \mathbf{x})$
- Use the identity to reweigh each data example:

$$\int \int l(f(\mathbf{x}), y) p(y \mid \mathbf{x}) p(\mathbf{x}) d\mathbf{x} dy = \int \int l(f(\mathbf{x}), y)$$

- **Weight Calculation:**

- Reweigh each example by the ratio .

- Train the model using *weighted empirical risk minimization*:

eq\_weighted-empirical-risk-min

$$\underset{f}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \beta_i l(f(\mathbf{x}_i), y_i).$$

- **Estimating the Ratio:**

- Estimate the ratio .  
 $\frac{p(\mathbf{x})}{q(\mathbf{x})}$
- Requires samples from both distributions: true (e.g., test data) and training set .
- Only need features , not labels .  
 $p(\mathbf{x})$   
 $q(\mathbf{x})$

- **Effective Approach:** Logistic Regression

- Learn a classifier to distinguish between data drawn from and data drawn from .
- If instances are indistinguishable, they are equally likely from either distribution.  
 $p(\mathbf{x})$   
 $q(\mathbf{x})$
- Well-discriminated instances should be significantly overweighted or underweighted accordingly.

# Covariate Shift Correction

- **Assumption:** Equal number of instances from distributions  $p$  and  $q$ .
- **Labels:**

- for data from  $p$ .
- $\bar{z} = 1$  for data from  $q$

- **Probability in Mixed Dataset:**

- 
- 

$$P(z = 1 \mid \mathbf{x}) = \frac{p(\mathbf{x})}{p(\mathbf{x}) + q(\mathbf{x})}$$

- **Logistic Regression Approach:**

- 
- Weight Calculation:

$$P(z = 1 \mid \mathbf{x}) = \frac{1}{1 + \exp(-h(\mathbf{x}))}$$

$$\beta_i = \exp(h(\mathbf{x}_i))$$



- **Steps to Correct Covariate Shift:**
  1. Distinguish between data from both distributions.
  2. Perform weighted empirical risk minimization using weights . $\beta_i$
- **Correction Algorithm:**
  - **Training Set:**
  - **Unlabeled Test Set:**  $\{(x_1, y_1), \dots, (x_n, y_n)\}$
  - **Assumption:**  $\{u_1, \dots, u_m\}$ 
    - $x_i$  drawn from source distribution.
    - $x_i$  for  $1 \leq i \leq n$  drawn from target distribution.
  - **Prototypical Algorithm:**
    1. **Step 1:** Use logistic regression to estimate the weights .
    2. **Step 2:** Apply weighted empirical risk minimization with the weights . $\beta_i$

1. Generate a binary-classification training set: .
  2. Train a binary classifier using logistic regression to get function  $\{(x_1, -1), (x_n, -1), (u_1, 1), \dots, (u_m, 1)\}$
  3. Weigh training data using  $\beta_i$  or better for some constant .  
$$\beta_i = \frac{1}{h(x_i)} \ln \left( \frac{\exp(h(x_i))}{\exp(h(x_i)) + c} \right)$$
  4. Use weights  $\beta_i$  for training on  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  in  $\text{eq_weighted-empirical-risk-min.}$
- Crucial Assumption:** - For this scheme to work, each data example in the target distribution must have a nonzero probability of occurring at training time. - If  $p(\mathbf{x})q(\mathbf{x}) = 0$ , the corresponding importance weight would be infinity.

# Label Shift Correction

- **Assumption:** Classification task with  $k$  categories.

- and are the source (training) and target (test) distributions, respectively.
  - Label distribution shifts: .
  - Class-conditional distribution remains the same: .
$$q(\mathbf{x} \mid y) = p(\mathbf{x} \mid y)$$
- **Risk Correction:**
  - Use the following identity:

$$\int \int l(f(\mathbf{x}), y) p(\mathbf{x} \mid y) p(y) \, d\mathbf{x} dy = \int \int l(f(\mathbf{x}), y) q(\mathbf{x} \mid y) q(y) \, d\mathbf{x} dy$$

- **Label Likelihood Ratios:**

- Importance weights: .

- **Estimation of Weights:**  $\beta_i = \frac{p(y_i)}{q(y_i)}$

- Use a good model on the source distribution to get consistent estimates of these weights.
- Inputs are high-dimensional (e.g., images), labels are simpler (e.g., categories).

- **Confusion Matrix:**

- Compute confusion matrix using the validation set.
- is a matrix: columns are true labels, rows are predicted categories.
- Each cell  $C_{ij}$  is the fraction of predictions where the true label was  $i$  and predicted label was  $j$ .

- **Estimating Target Label Distribution:**

- Average model predictions on test data to get mean model outputs .

- : fraction of predictions where the model predicted .  
 $\mu(\hat{y}_i)$   
 $i$

- **Linear System:**

- Estimate test set label distribution by solving .

$$\mathbf{C}p(\mathbf{y}) = \mu(\hat{\mathbf{y}})$$

- Solution: if  $\mathbf{C}$  is invertible.
- **Calculating Weights:**  $\mathbf{C}^{-1} \bar{\mathbf{C}}^T \mu(\hat{\mathbf{y}})$

- Estimate from source data.

- For any training example  $i$  with label  $y_i$ , calculate weight .

- Use weights in weighted empirical risk minimization in  $\beta_i$ .  

$$\min_{\beta} \sum_i \bar{w}_i p(y_i)$$



# Concept Shift Correction

- **Challenge:**
  - Concept shift is harder to correct in a principled manner.
  - Example: Shift from distinguishing cats vs. dogs to white vs. black animals may require new labels and training from scratch.
  - Extreme shifts are rare; usually, the task changes slowly over time.
- **Examples:**
  - **Computational Advertising:** New products are launched, old products become less popular. Distribution over ads and their popularity changes gradually.
  - **Traffic Cameras:** Lenses degrade gradually due to environmental wear, affecting image quality progressively.
  - **News Content:** Changes gradually with new stories appearing while most news remains unchanged.

- **Approach:**

- Use existing network weights and perform a few update steps with new data to adapt to gradual changes, rather than training from scratch.

# A Taxonomy of Learning Problems

- **Overview:**
  - Consider other aspects of machine learning problem formulation with knowledge on dealing with distribution changes.

## Batch Learning

- **Definition:**
  - Training a model using a batch of training features and labels .
  - Model is deployed to score new data drawn from the same distribution.  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
- **Example:**
  - Train a cat detector with pictures of cats and dogs.
  - Deploy as part of a smart catdoor system, which is then installed and never updated again.

## Online Learning

- **Definition:**
  - Data arrives one sample at a time.
  - Observe  $(\mathbf{x}_i, y_i)$ , estimate , then observe and update the model based on the incurred loss.

- **Example:**

- Predict tomorrow's stock price, trade based on the estimate, and learn from the profit or loss at the end of the day.

- **Cycle:**

model  $f_t \rightarrow$  data  $\mathbf{x}_t \rightarrow$  estimate  $f_t(\mathbf{x}_t) \rightarrow$   
observation  $y_t \rightarrow$  loss  $l(y_t, f_t(\mathbf{x}_t)) \rightarrow$  model  $f_{t+1}$

## Control

- **Definition:**
  - Environment remembers previous actions, influencing current observations.
- **Examples:**
  - **PID Controller:** Adjusts heating based on previous boiler temperatures.
  - **User Behavior:** Depends on previously shown news articles.
- **Applications:**
  - Model the environment to make decisions less random.
  - Used for automatic hyperparameter tuning to improve model performance  
:cite:**Shao.Yao.Sun.ea.2020**.

## Considering the Environment

- **Key Distinction:**
  - Strategies for stationary environments may not work when environments adapt.
- **Example:**
  - Arbitrage opportunity disappears once exploited by a trader.
- **Algorithm Selection:**
  - Environment changes determine suitable algorithms.
  - Slow changes: Force estimates to change slowly.
  - Instantaneous but infrequent changes: Make allowances for sudden shifts.
- **Crucial Knowledge:**
  - Understanding environment dynamics is essential for dealing with concept shift.

# Summary

- **Distribution Shift:**
  - Training and test sets may come from different distributions.
  - This mismatch is known as distribution shift.
- **Risk and Empirical Risk:**
  - Risk: Expectation of the loss over the entire population drawn from the true distribution (usually unavailable).
  - Empirical Risk: Average loss over the training data, used to approximate risk.
  - Practical approach: Perform empirical risk minimization.

- **Covariate and Label Shift:**
  - These types of shifts can be detected and corrected for at test time under corresponding assumptions.
  - Ignoring these biases can lead to problematic performance during testing.
- **Environment Response:**
  - The environment may remember automated actions and respond in unexpected ways.
  - Important to consider this when building models.
  - Continuously monitor live systems to ensure models and environments do not become entangled in unanticipated ways.

