

Bayesian Learning of Classifiers with Stan

Torsten Scholak

twitter [@tscholak](#) | github [tscholak](#) | blog [tscholak.github.io](#)

March 9, 2016

<http://x.x.x.x:8000>

Outline

Classification

Discriminative Approach

Generative Approach

Bayesian Inference in Stan

Demo

<http://x.x.x.x:8000>

The Classification Problem

Have

D -dimensional, real explanatory variable $\mathbf{X} = (X_1, \dots, X_D)^\top$

Class label C , assumes one out of K classes $\{1, \dots, K\}$

Training set of N observations $(\mathbf{x}, c)_{1:N}$

Want

Assign each new observation $\mathbf{x}'_{n'}, n' = 1, \dots, N'$, to *one* class

$$\mathbf{x}'_{n'} \rightarrow \mathbf{x}'_{n'}, c'_{n'}$$

How can we infer the labels $c'_{1:N'}$?

Discriminative Approach

Train a model so as to maximize the probability of getting correct labels

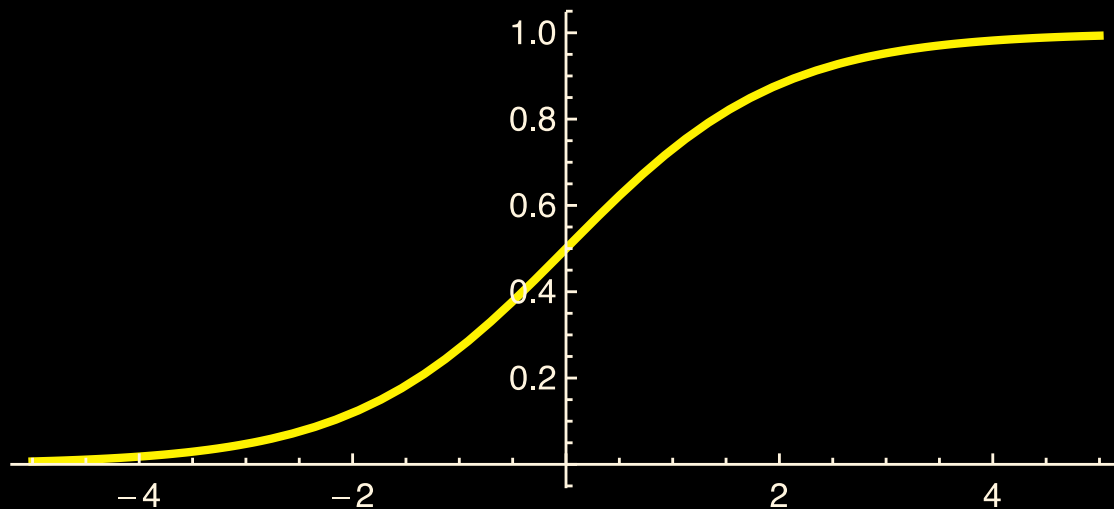
Logistic "Regression"

Logistic regression estimates $p(c = k | \mathbf{x})$, $k = 1, \dots, K$

Idea: Map the output of a linear model, $b + \mathbf{w}^\top \mathbf{x}$, to a probability

$$p(c = k | \mathbf{x}, \mathbf{b}, \mathbf{W}) = \text{softmax}_k(\mathbf{b} + \mathbf{W}\mathbf{x}) = \frac{\exp(b_k + \mathbf{w}_k^\top \mathbf{x})}{\sum_{k'=1}^K \exp(b_{k'} + \mathbf{w}_{k'}^\top \mathbf{x})}$$

Softmax is a soft activation function with probability interpretation:



Graph of $y = (1 + \exp(-x))^{-1}$

Training

Maximum Likelihood Estimation

Maximize the conditional log-likelihood of the training data $(\mathbf{x}, c)_{1:N}$:

$$\operatorname{argmax}_{\mathbf{b}, \mathbf{W}} \sum_{n=1}^N \log p(c_n | \mathbf{x}_n, \mathbf{b}, \mathbf{W})$$

Maximum A Posteriori Estimation

Add ℓ_2 -regularization and maximize:

$$\operatorname{argmax}_{\mathbf{b}, \mathbf{W}} \sum_{n=1}^N \log p(c_n | \mathbf{x}_n, \mathbf{b}, \mathbf{W}) - \lambda \sum_{k=1}^K (|b_k|^2 + \|\mathbf{w}_k\|^2) =$$

$$\operatorname{argmax}_{\mathbf{b}, \mathbf{W}} \prod_{n=1}^N p(c_n | \mathbf{x}_n, \mathbf{b}, \mathbf{W}) \times \prod_{k=1}^K [\text{Normal}(b_k | 0, \frac{1}{2\lambda}) \text{Normal}(\mathbf{w}_k | \mathbf{0}, \frac{1}{2\lambda})]$$

Prediction

Discrimination is based on hard boundaries:

$$\begin{aligned} c'_{n'} &= \operatorname{argmax}_k p(c'_{n'} = k \mid \mathbf{x}'_{n'}, \mathbf{b}, \mathbf{W}) \\ &= \operatorname{argmax}_k \operatorname{softmax}_k(\mathbf{b} + \mathbf{W}\mathbf{x}'_{n'}) \end{aligned}$$

Problems

How to choose the regularization parameter λ so as to avoid overfitting?

How to incorporate prior knowledge?

How to integrate this method into a larger model?

Isn't discrimination hiding the fact that there are uncertainties in the prediction?

Bayesian Generative Approach

Build a full probabilistic model of all variables, not only the class label

Bayesian Learning

1. Introduce a *generative model* of the data, $\mathbf{x}_{1:N}$, conditioned on the class labels $c_{1:N}$ and other latent variables θ :

$$p(\mathbf{x}_{1:N} | c_{1:N}, \theta)$$

2. Model a *prior probability* for the latent variables:

$$p(c_{1:N}, \theta)$$

3. Together, that will give us the *posterior probability*:

$$p(c_{1:N}, \theta | \mathbf{x}_{1:N}) \propto p(\mathbf{x}_{1:N} | c_{1:N}, \theta) p(c_{1:N}, \theta)$$

(Bayes theorem)

Why Bayesian Learning?

Resistant to noise, avoids overfitting

Takes into account prior knowledge, better results in smaller samples

More flexibility, straightforward integration into larger model

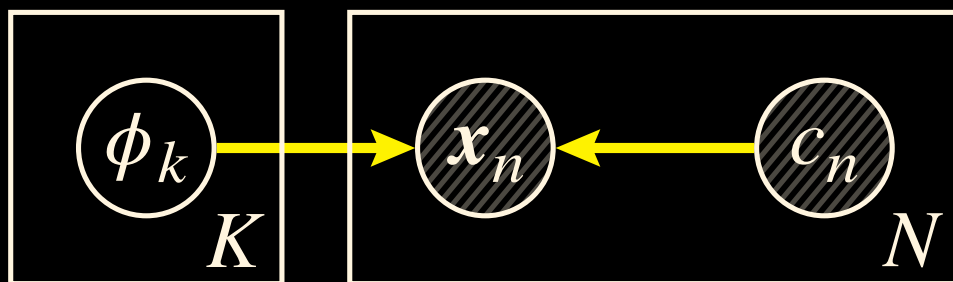
Predictions are probabilistic rather than discriminative

A General Generative Model for Classification

We imagine that

$$\mathbf{x}_n \mid c_n, \phi_{1:K} \sim f(\phi_{c_n})$$

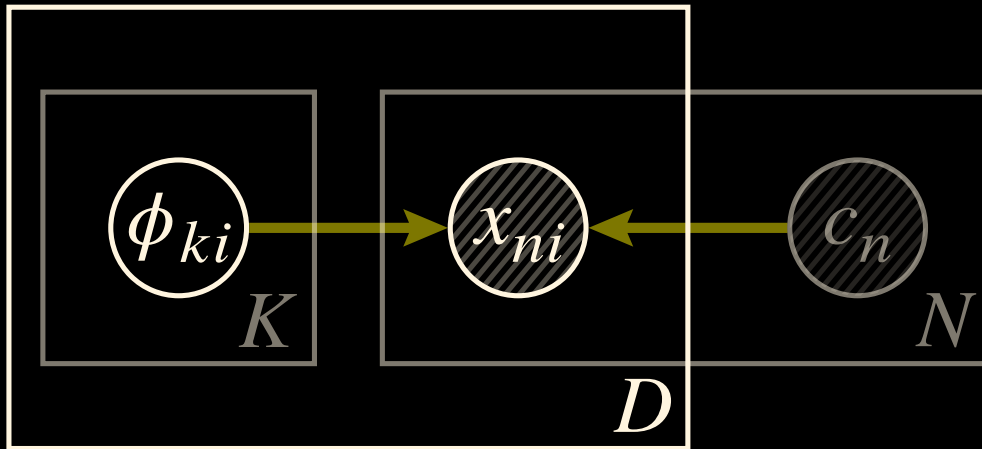
In words: *Given the label c_n , each observed data point \mathbf{x}_n is drawn from some distribution with probability density $f(\phi_{c_n})$, where the ϕ_k are latent parameters*



A Naïve Assumption

As a simplification, assume that the features factorize:

$$p(\mathbf{x}_n | c_n, \phi_{1:K}) = \prod_{i=1}^D f(x_{ni} | \phi_{c_n i})$$

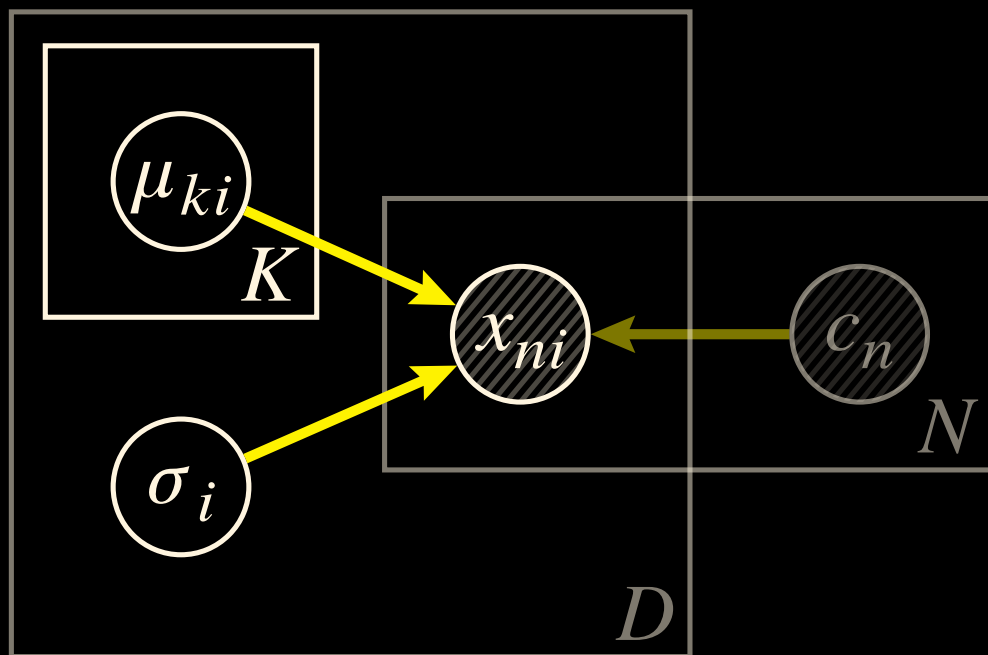


Generate Normally Distributed Samples

For reasons which will become clear soon, choose a normal distribution:

$$\begin{aligned} f(x_i | \phi_{ki}) &\equiv \text{Normal}(x_i | \mu_{ki}, \sigma_i^2) \\ &= \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left[-\frac{x_i - \mu_{ki}}{\sigma_i}\right] \end{aligned}$$

with latent means μ_{ki} and variances σ_i^2

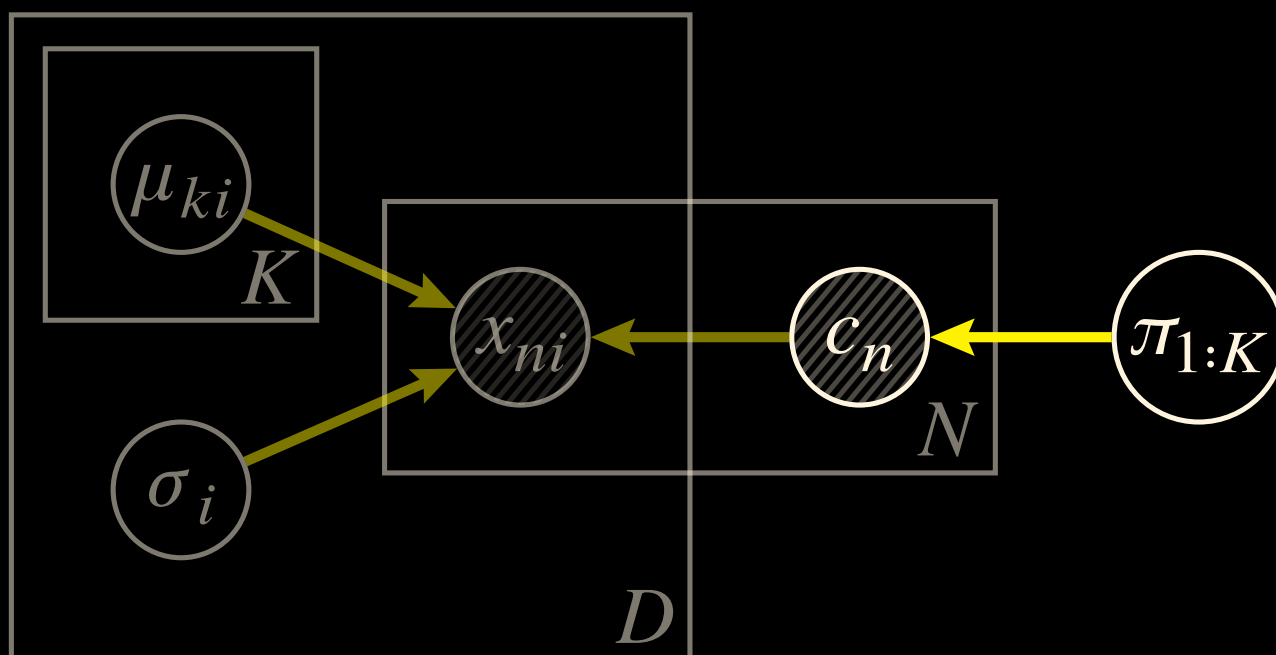


A Class Indicator Prior

We further imagine that

$$c_n \mid \pi_{1:K} \sim \text{Categorical}(\pi_{1:K})$$

In words: *The indicator variables are distributed according to a categorical distribution defined on event probabilities π_k with $0 \leq \pi_k \leq 1$ and $\sum_{k=1}^K \pi_k = 1$*

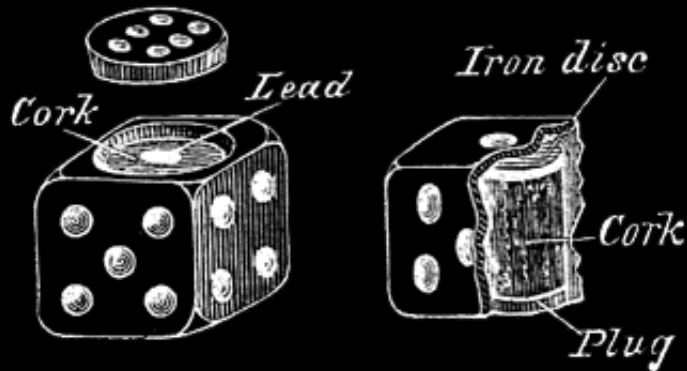


The Categorical Distribution

Probability Mass Function

Loaded K -sided die roll:

$$p(c_n = k | \pi_{1:K}) = \pi_k$$



The Connection to Logistic Regression

With a naïve Gaussian prior on \mathbf{X} and a categorical prior on C , we have:

$$p(c_n = k | \mathbf{x}_n, \pi_{1:K}, \boldsymbol{\mu}_{1:K}, \boldsymbol{\sigma}) = \frac{\pi_k \text{Normal}(\mathbf{x}_n | \boldsymbol{\mu}_k, \text{diag}(\boldsymbol{\sigma}^2))}{\prod_{k'=1}^K \pi_{k'} \text{Normal}(\mathbf{x}_n | \boldsymbol{\mu}_{k'}, \text{diag}(\boldsymbol{\sigma}^2))}$$

$$\vdots$$

$$= \text{softmax}_k(\mathbf{b} + \mathbf{W}\mathbf{x}_n)$$

with

$$b_k = \log \pi_k - \sum_{i=1}^D \frac{\mu_{ki}^2}{2\sigma_i^2}$$

$$w_{ki} = \frac{\mu_{ki}}{\sigma_i^2}$$

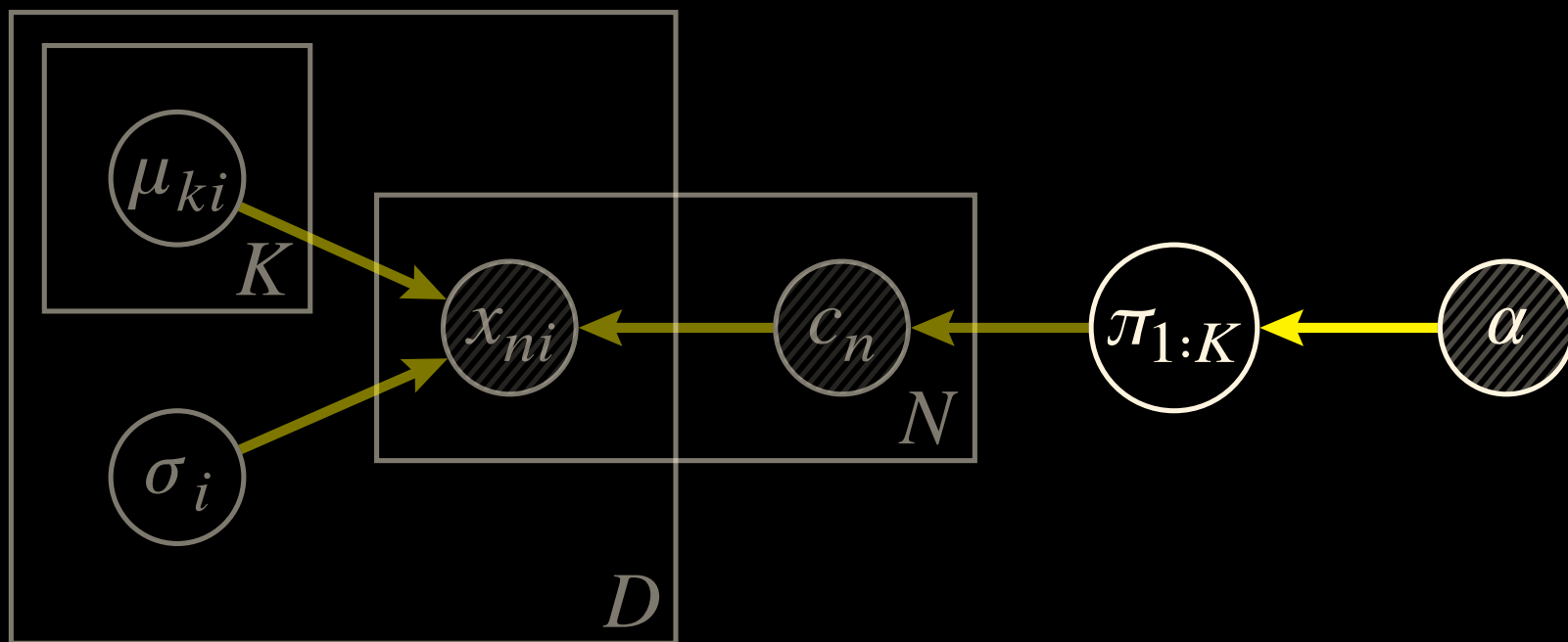
That's the well-known likelihood function of logistic regression!

A Class Prevalence Prior

We choose a *conjugate prior* for the probabilities:

$$\pi_{1:K} \mid \alpha \sim \text{Dirichlet}(\alpha)$$

In words: *The probabilities π_k are distributed according to a symmetric Dirichlet distribution with concentration parameter $\alpha > 0$*

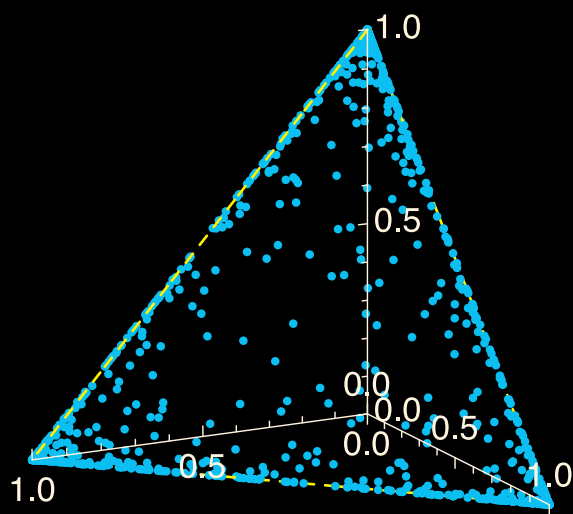


The Dirichlet Distribution I

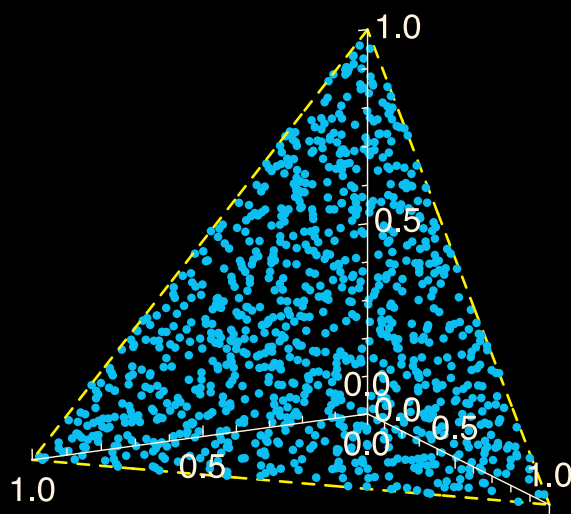
Probability Density Function

For the symmetric Dirichlet distribution:

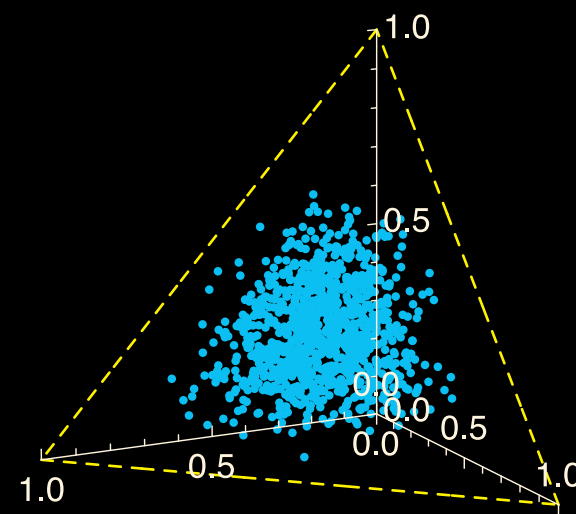
$$p(\pi_{1:K}|\alpha) = \frac{\Gamma(K\alpha)}{\Gamma(\alpha)^K} \prod_{k=1}^K \pi_k^{\alpha-1}$$



If $\alpha \ll 1$, then concentrated around the corners of the simplex



If $\alpha = 1$, then uniformly distributed over the simplex



If $\alpha \gg 1$, then concentrated around the center of the simplex

The Dirichlet Distribution II

Pólya's urn

Consider an urn containing balls of K different colors.

Initially, the urn contains each α balls of colors $1, 2, \dots, K$.

Now perform N draws from the urn, where after each draw, the ball is placed back into the urn with an additional ball of the same color.

In the limit as $N \rightarrow \infty$, the proportions $\pi_{1:K}$ of different colored balls in the urn will be distributed as $\text{Dirichlet}(\alpha)$.

Pick distributions for the μ_{ki} and σ_i :

$$\sigma_i^{-2} | a_0, b_0 \sim \text{Gamma}(a_0, b_0)$$

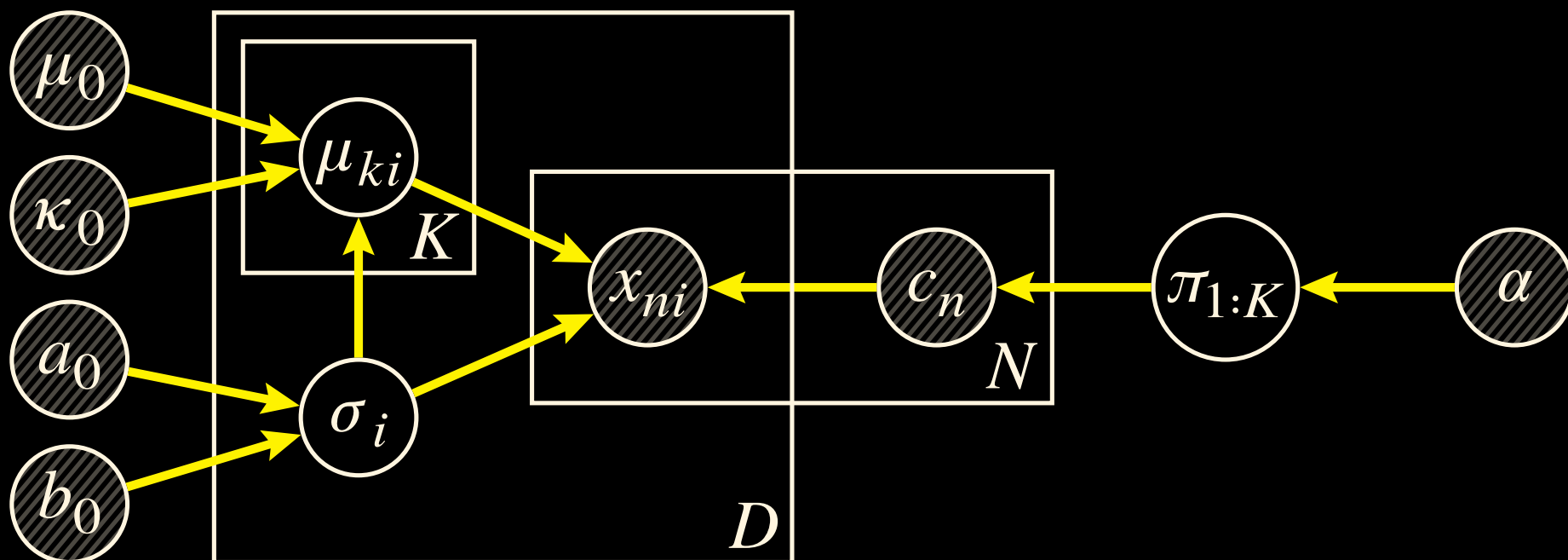
Figure 1: A graphical model for the proposed method. The model is divided into three main sections: a latent variable section (D), a data section (N), and a hyperparameter section (K). The latent variable section (D) contains nodes μ_{ki} and σ_i . The data section (N) contains nodes x_{ni} and c_n . The hyperparameter section (K) contains nodes μ_0 , κ_0 , a_0 , b_0 , $\pi_{1:K}$, and α . The nodes are connected by directed edges: μ_0 and κ_0 point to μ_{ki} ; a_0 and b_0 point to σ_i ; μ_{ki} points to x_{ni} ; σ_i points to x_{ni} ; c_n points to x_{ni} ; $\pi_{1:K}$ points to c_n ; and α points to $\pi_{1:K}$. The nodes are shaded gray, and the edges are blue.

The Posterior of The Complete Model

$$p(c_{1:N}, \pi_{1:K}, \mu_{1:K}, \sigma \mid x_{1:N}, \alpha, \mu_0, \kappa_0, a_0, b_0)$$

$$\propto \prod_{n=1}^N \left[\prod_{i=1}^D p(x_{ni} \mid \mu_{c_n i}, \sigma_i^2) \right] p(c_n \mid \pi_{1:K})$$

$$\times p(\pi_{1:K} \mid \alpha) \prod_{i=1}^D \left[\prod_{k=1}^K p(\mu_{ki} \mid \mu_0, \kappa_0) \right] p(\sigma_i^{-2} \mid a_0, b_0)$$



Bayesian Inference in Stan

Building, tweaking, enhancing, and enriching models easily without having to think about the implementation

Stan

1. Imperative probabilistic programming language
2. Automatic differentiation for HMC
3. Adaptation routines
4. R, Python, MATLAB, Julia, Stata, and command line interfaces

Program Blocks of A Stan Program

Declare data and parameter variables, define the log-posterior:

```
data { }  
transformed data { }  
parameters { }  
transformed parameters { }  
model { }  
generated quantities { }
```

The data Block

Reading in information from an external source:

```
data {  
  int<lower=1> K;           // number of classes  
  int<lower=1> D;           // number of features  
  int<lower=0> N;           // number of labelled observations  
  int<lower=1,upper=K> c[N]; // classes for labelled observations  
  vector[D] x[N];          // features for labelled observations  
  vector<lower=0>[K] alpha; // class concentration  
  real mu0;                 // prior mean  
  real<lower=0> kappa0;      // prior sample size  
  real<lower=0> a0;          // shape  
  real<lower=0> b0;          // rate  
}
```

The transformed data Block

Manipulate the external information once:

```
transformed data { }
```

We won't need it

The parameters Block

Define the things we are going to sample from:

```
parameters {  
    simplex[K] pi;           // class prevalence  
    vector[D] mu[K];         // means of features  
    vector<lower=0>[D] invsigmasqr; // inverse variances of features  
}
```

The transformed parameters Block

Process parameters before computing the posterior:

```
transformed parameters {  
  vector<lower=0>[D] sigma;    // scales of features  
  for (i in 1:D)  
    sigma[i] <- inv_sqrt(invsigmasqr[i]);  
}
```

Parameters defined here are not sampled by the Markov chain

The model Block

Define the posterior:

```
model {  
  pi ~ dirichlet(alpha);      // class prevalence prior  
  for (n in 1:N)  
    c[n] ~ categorical(pi);  // class indicator prior  
  for (k in 1:K)  
    mu[k] ~ normal(mu0, sigma/sqrt(kappa0));  
  for (i in 1:D)  
    invsigmasqr[i] ~ gamma(a0, b0);  
                                // base measures  
  for (n in 1:N)  
    x[n] ~ normal(mu[c[n]], sigma);  
                                // generative model  
}
```

The generated quantities Block

Produce random samples, e.g. to validate the model with pseudo-data:

```
generated quantities {  
  vector[K] z[Np];  
  vector[K] sm[Np];  
  int<lower=1,upper=K> cp[Np];  
  for (np in 1:Np) {  
    for (k in 1:K)  
      z[np, k] <- normal_log(xp[np], mu[k], sigma) \  
        + categorical_log(k, pi);  
    sm[np] <- softmax(z[np]);  
    cp[np] <- categorical_rng(sm[np]);  
  }  
}
```


Demo

Further Content Ingestion

[Stan User Guide and Reference Manual v2.9.0](#)

Tons of example programs

Introduction to Bayesian statistics

Language reference

Discussion of HMC and NUTS algorithms

[Example program repository](#)

[Michael Betancourt's YouTube videos](#)

[Andrew Gelman's book](#)