

Assignment 2: Stochastic Variational Inference in the TrueSkill Model

STA414/STA2014 and CSC412/CSC2506 Winter 2020

Justin Leo

Student Number: 1006376459

March 18, 2020

The goal of this assignment is to get you familiar with the basics of Bayesian inference in large models with continuous latent variables, and the basics of stochastic variational inference.

Background We'll implement a variant of the TrueSkill model, a player ranking system for competitive games originally developed for Halo 2. It is a generalization of the Elo rating system in Chess. For the curious, the original 2007 NIPS paper introducing the trueskill paper can be found here: <http://papers.nips.cc/paper/3079-trueskilltm-a-bayesian-skill-rating-system.pdf>

This assignment is based on one developed by Carl Rasmussen at Cambridge for his course on probabilistic machine learning: <http://mlg.eng.cam.ac.uk/teaching/4f13/1920/>

0.1 Model definition

We'll consider a slightly simplified version of the original trueskill model. We assume that each player has a true, but unknown skill $z_i \in \mathbb{R}$. We use N to denote the number of players.

The prior. The prior over each player's skill is a standard normal distribution, and all player's skills are *a priori* independent.

The likelihood. For each observed game, the probability that player i beats player j , given the player's skills z_A and z_B , is:

$$p(A \text{ beat } B | z_A, z_B) = \sigma(z_i - z_j)$$

where

$$\sigma(y) = \frac{1}{1 + \exp(-y)}$$

There can be more than one game played between a pair of players, and in this case the outcome of each game is independent given the players' skills. We use M to denote the number of games.

The data. The data will be an array of game outcomes. Each row contains a pair of player indices. The first index in each pair is the winner of the game, the second index is the loser. If there were M games played, then the array has shape $M \times 2$.

1 Implementing the model [10 points]

- (a) [2 points] Implement a function `log_prior` that computes the log of the prior over all player's skills. Specifically, given a $K \times N$ array where each row is a setting of the skills for all N players, it returns a $K \times 1$ array, where each row contains a scalar giving the log-prior for that set of skills.

```
function log_prior(zs)
    prior = factorized_gaussian_log_density(0,0,zs)
    return prior
end
```

- (b) [3 points] Implement a function `logp_a_beats_b` that, given a pair of skills z_a and z_b evaluates the log-likelihood that player with skill z_a beat player with skill z_b under the model detailed above. To ensure numerical stability, use the function `log1pexp` that computes $\log(1 + \exp(x))$ in a numerically stable way. This function is provided by `StatsFuns.jl` and imported already, and also by Python's `numpy`.

```
function logp_a_beats_b(za,zb)
    x = zb-za
    return (-log1pexp.(x))
end
```

- (c) [3 points] Assuming all game outcomes are i.i.d. conditioned on all players' skills, implement a function `all_games_log_likelihood` that takes a batch of player skills `zs` and a collection of observed games `games` and gives a batch of log-likelihoods for those observations. Specifically, given a $K \times N$ array where each row is a setting of the skills for all N players, and an $M \times 2$ array of game outcomes, it returns a $K \times 1$ array, where each row contains a scalar giving the log-likelihood of all games for that set of skills. Hint: You should be able to write this function without using for loops, although you might want to start that way to make sure what you've written is correct. If A is an array of integers, you can index the corresponding entries of another matrix B for every entry in A by writing $B[A]$.

```
function all_games_log_likelihood(zs,games)
    m = logp_a_beats_b(zs[games[:,1],:], zs[games[:,2],:])
    return sum(m, dims=1)
end
```

- (d) [2 points] Implement a function `joint_log_density` which combines the log-prior and log-likelihood of the observations to give $p(z_1, z_2, \dots, z_N, \text{all game outcomes})$

```
function joint_log_density(zs,games)
    return log_prior(zs) .+ all_games_log_likelihood(zs,games)
end
```

```
@testset "Test shapes of batches for likelihoods" begin
    B = 15 # number of elements in batch
    N = 4 # Total Number of Players
    test_zs = randn(4,15)
    test_games = [1 2; 3 1; 4 2] # 1 beat 2, 3 beat 1, 4 beat 2
    @test size(test_zs) == (N,B)
    #batch of priors
    @test size(log_prior(test_zs)) == (1,B)
    # loglikelihood of p1 beat p2 for first sample in batch
    @test size(logp_a_beats_b(test_zs[1,1],test_zs[2,1])) == ()
    # loglikelihood of p1 beat p2 broadcasted over whole batch
    @test size(logp_a_beats_b.(test_zs[1,:],test_zs[2,:])) == (B,)
    # batch loglikelihood for evidence
```

```

@test size(all_games_log_likelihood(test_zs,test_games)) == (1,B) #B changed to N
# batch loglikelihood under joint of evidence and prior
@test size(joint_log_density(test_zs,test_games)) == (1,B)
end

```

Test Summary:	Pass	Total
Test shapes of batches for likelihoods	6	6

2 Examining the posterior for only two players and toy data [10 points]

To get a feel for this model, we'll first consider the case where we only have 2 players, A and B . We'll examine how the prior and likelihood interact when conditioning on different sets of games.

Provided in the starter code is a function `skillcontour!` which evaluates a provided function on a grid of z_A and z_B 's and plots the isocontours of that function. As well there is a function `plot_line_equal_skill!`. We have included an example for how you can use these functions.

We also provided a function `two_player_toy_games` which produces toy data for two players. I.e. `two_player_toy_games(5,3)` produces a dataset where player A wins 5 games and player B wins 3 games.

- (a) [2 points] For two players A and B , plot the isocontours of the joint prior over their skills. Also plot the line of equal skill, $z_A = z_B$. Hint: you've already implemented the **log** of the likelihood function.

```

# Example for how to use contour plotting code
plot(title="Log Prior Contour Plot",
      xlabel = "Player 1 Skill",
      ylabel = "Player 2 Skill", margin=5mm
    )
example_gaussian(zs) = exp(factorized_gaussian_log_density([-1.,2.],[0.,0.5],zs))
#skillcontour!(example_gaussian; label="example gaussian")
skillcontour!(example_gaussian)
plot_line_equal_skill!()

```

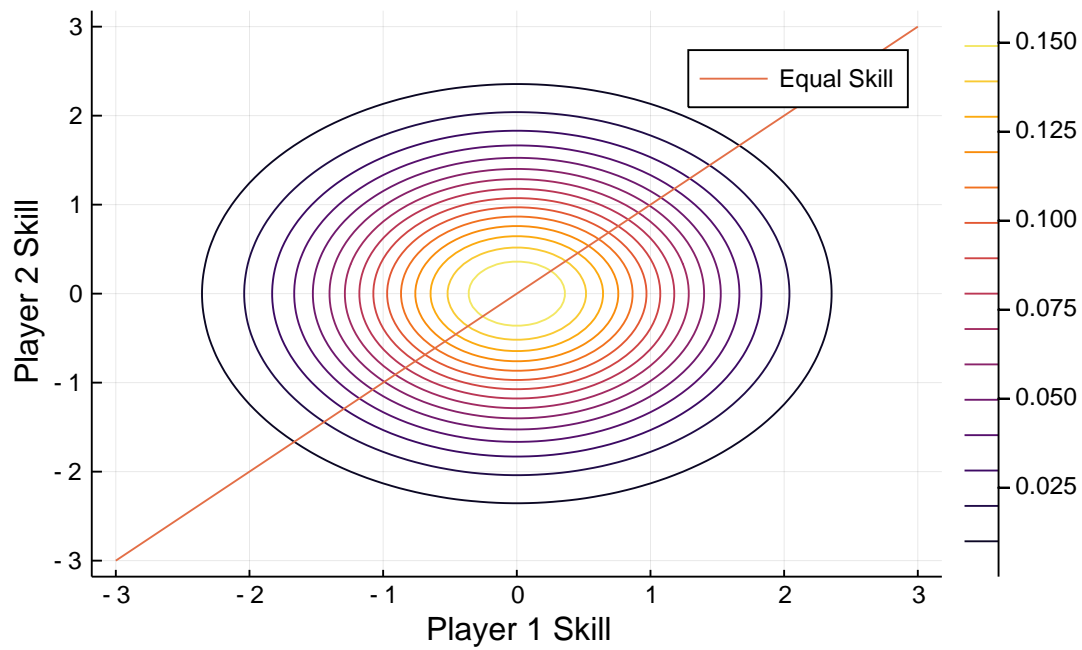
```

#prior contour:
#a

f1(zs) = exp.(log_prior(zs))
t = display(plot(title="Joint Prior Contour Plot",
                  xlabel = "Player 1 Skill",
                  ylabel = "Player 2 Skill", margin=5mm))
display(skillcontour!(f1))
display(plot_line_equal_skill!())

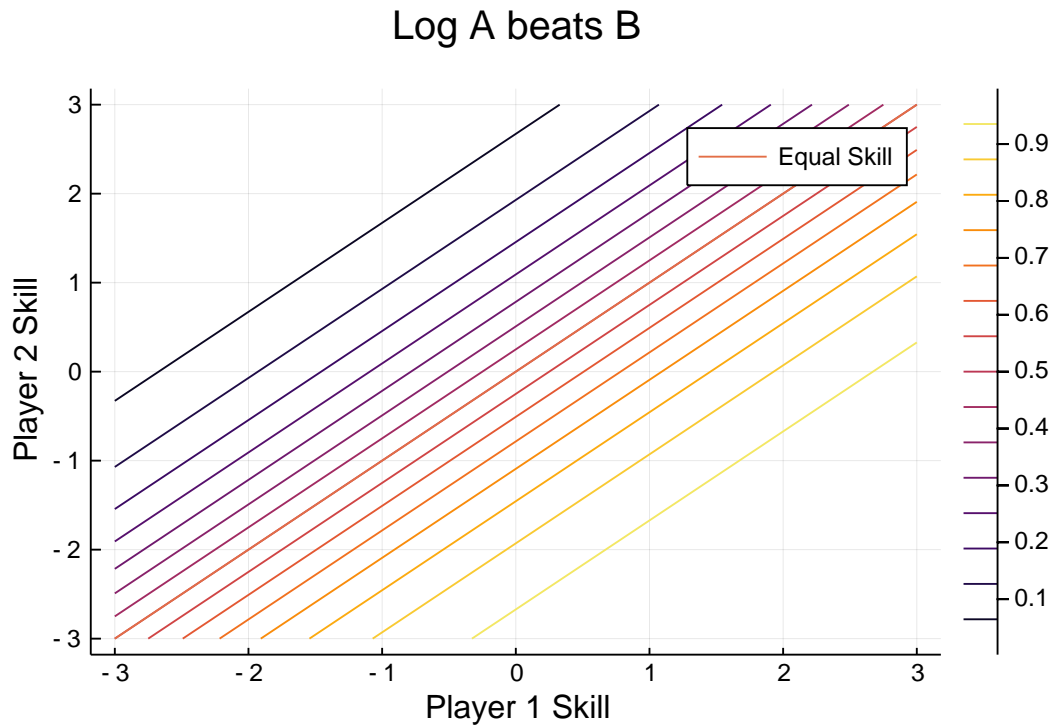
```

Joint Prior Contour Plot



(b) [2 points] Plot the isocontours of the likelihood function. Also plot the line of equal skill, $z_A = z_B$.

```
f11(zs) = exp.(logp_a_beats_b(zs[1,:],zs[2,:]))
t = display(plot(title="Log A beats B",
    xlabel = "Player 1 Skill",
    ylabel = "Player 2 Skill", margin=5mm))
    #prior(zs) = exp.(log_prior(zs))
display(skillcontour!(f11))
display(plot_line_equal_skill!())
```

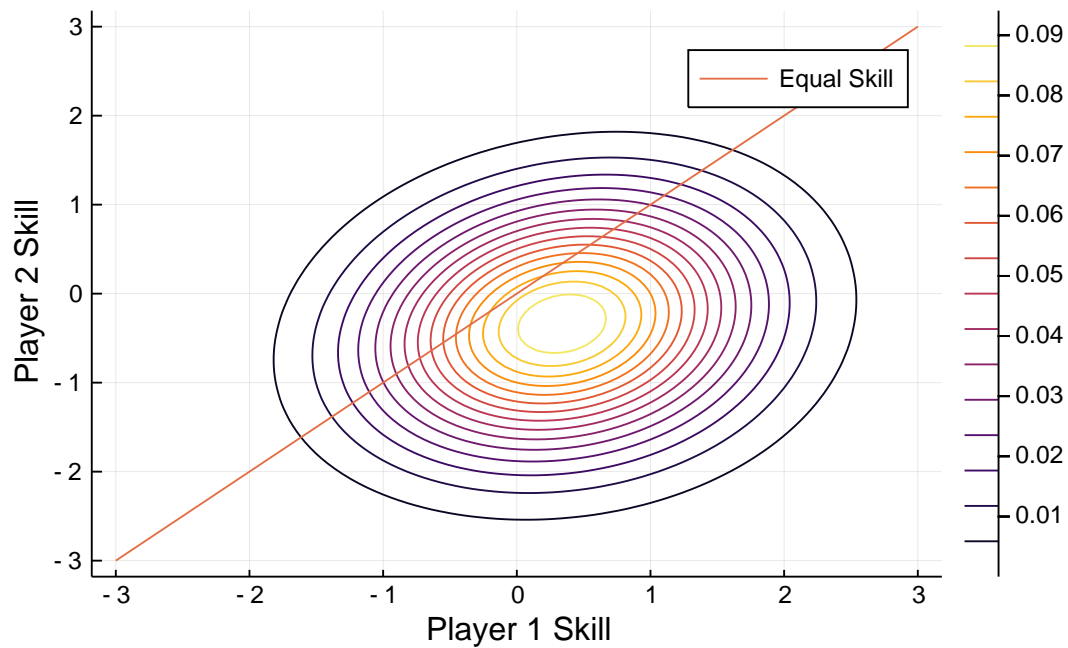


- (c) [2 points] Plot isocontours of the joint posterior over z_A and z_B given that player A beat player B in one match. Since the contours don't depend on the normalization constant, you can simply plot the isocontours of the log of joint distribution of $p(z_A, z_B, \text{A beat B})$. Also plot the line of equal skill, $z_A = z_B$.

```
a_win_1 = two_player_toy_games(1, 0)
```

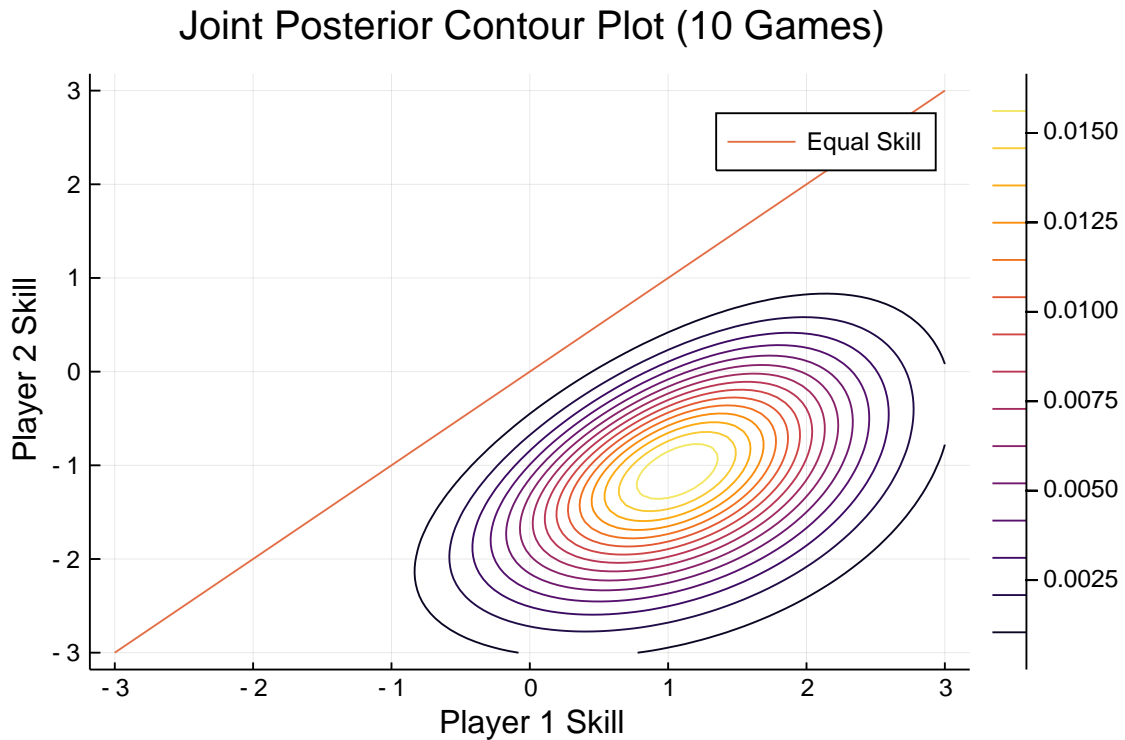
```
f2(zs) = exp(joint_log_density(zs, a_win_1))
t = display(plot(title="Joint Posterior Contour Plot (1 Game)",
  xlabel = "Player 1 Skill",
  ylabel = "Player 2 Skill", margin=5mm))
  #prior(zs) = exp.(log_prior(zs))
  display(skillcontour!(f2))
  display(plot_line_equal_skill!())
```

Joint Posterior Contour Plot (1 Game)



- (d) [2 points] Plot isocountours of the joint posterior over z_A and z_B given that 10 matches were played, and player A beat player B all 10 times. Also plot the line of equal skill, $z_A = z_B$.

```
a_win_10 = two_player_toy_games(10, 0)
f3(zs) = exp(joint_log_density(zs,a_win_10))
t = display(plot(title="Joint Posterior Contour Plot (10 Games)",
    xlabel = "Player 1 Skill",
    ylabel = "Player 2 Skill", margin=3mm))
    #prior(zs) = exp.(log_prior(zs))
display(skillcontour!(f3))
display(plot_line_equal_skill!())
```

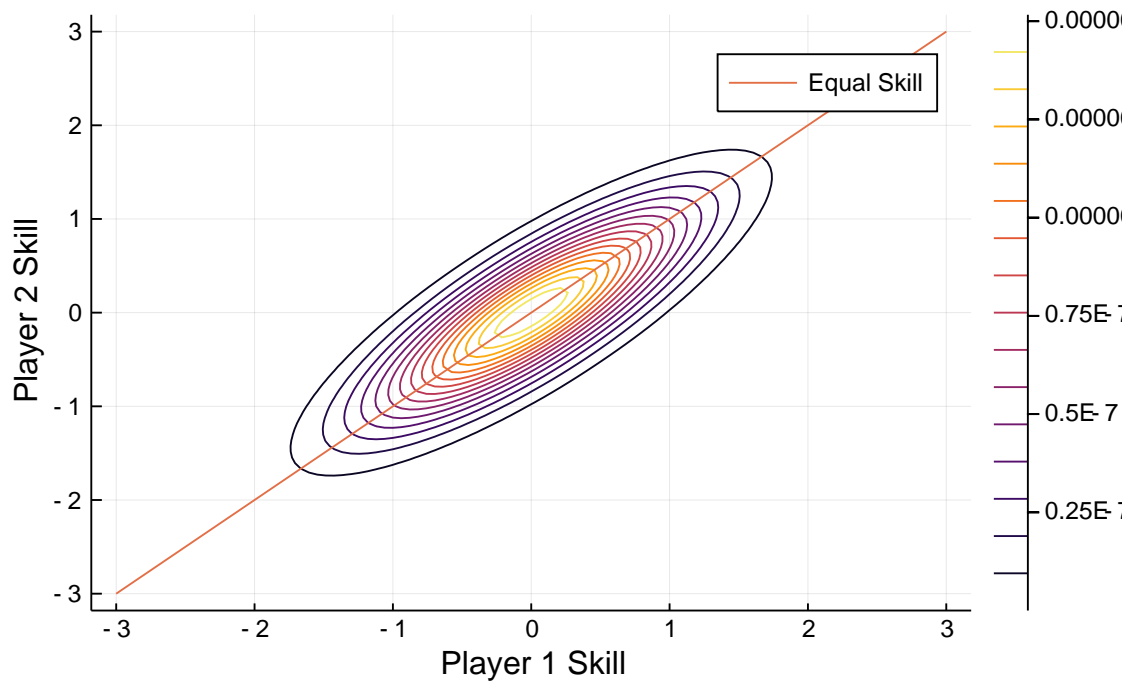


- (e) [2 points] Plot isocountours of the joint posterior over z_A and z_B given that 20 matches were played, and each player beat the other 10 times. Also plot the line of equal skill, $z_A = z_B$.

For all plots, label both axes.

```
split_match_10 = two_player_toy_games(10, 10)
f4(zs) = exp(joint_log_density(zs,split_match_10))
t = display(plot(title="Joint Posterior Contour Plot (20 Games)",
    xlabel = "Player 1 Skill",
    ylabel = "Player 2 Skill", margin=3mm))
    #prior(zs) = exp.(log_prior(zs))
display(skillcontour!(f4))
display(plot_line_equal_skill!())
```

Joint Posterior Contour Plot (20 Games)



3 Stochastic Variational Inference on Two Players and Toy Data [18 points]

One nice thing about a Bayesian approach is that it separates the model specification from the approximate inference strategy. The original Trueskill paper from 2007 used message passing. Carl Rasmussen's assignment uses Gibbs sampling, a form of Markov Chain Monte Carlo. We'll use gradient-based stochastic variational inference, which wasn't invented until around 2014.

In this question we will optimize an approximate posterior distribution with stochastic variational inference to approximate the true posterior.

- (a) [5 points] Implement a function `elbo` which computes an unbiased estimate of the evidence lower bound. As discussed in class, the ELBO is equal to the KL divergence between the true posterior $p(z|\text{data})$, and an approximate posterior, $q_\phi(z|\text{data})$, plus an unknown constant. Use a fully-factorized Gaussian distribution for $q_\phi(z|\text{data})$. This estimator takes the following arguments:

- `params`, the parameters ϕ of the approximate posterior $q_\phi(z|\text{data})$.
- A function `logp`, which is equal to the true posterior plus a constant. This function must take a batch of samples of z . If we have N players, we can consider B -many samples from the joint over all players' skills. This batch of samples `zs` will be an array with dimensions (N, B) .
- `num_samples`, the number of samples to take.

This function should return a single scalar. Hint: You will need to use the reparameterization trick when sampling `zs`.

```
function elbo(params, logp, num_samples)
    x = randn(size(params[1])[1], num_samples)
    a = exp.(params[2])
    b = params[1]
    samples = a .* x .+ b
    logp_estimate = logp(samples)
    logq_estimate = factorized_gaussian_log_density(params[1], params[2], samples)
    return sum(logp_estimate - logq_estimate) / (num_samples)
end
```

- (b) [2 points] Write a loss function called `neg_toy_elbo` that takes variational distribution parameters and an array of game outcomes, and returns the negative elbo estimate with 100 samples.

```
function neg_toy_elbo(params; games = two_player_toy_games(1,0), num_samples = 100)
    logp(zs) = joint_log_density(zs, games)
    return -elbo(params, logp, num_samples)
end
```

- (c) [5 points] Write an optimization function called `fit_toy_variational_dist` which takes initial variational parameters, and the evidence. Inside it will perform a number of iterations of gradient descent where for each iteration :

- Compute the gradient of the loss with respect to the parameters using automatic differentiation.
- Update the parameters by taking an `lr`-scaled step in the direction of the descending gradient.
- Report the loss with the new parameters (using `@info` or `print` statements)
- On the same set of axes plot the target distribution in red and the variational approximation in blue.

Return the parameters resulting from training.

```

num_players_toy = 2
toy_mu = [-2.,3.] # Initial mu, can initialize randomly!
toy_ls = [0.5,0.] # Initial log_sigma, can initialize randomly!
toy_params_init = (toy_mu, toy_ls)

function fit_toy_variational_dist(init_params, toy_evidence; num_itrs=200, lr= 1e-2, num_q_samples = 10)
    params_cur = init_params
    for i in 1:num_itrs
        grad_params = gradient((params) -> neg_toy_elbo(params; games = toy_evidence, num_samples = num_q_samples), params_cur)[1]
        params_cur = params_cur .- lr .* grad_params
        @info "Loss $(neg_toy_elbo(params_cur; games = toy_evidence, num_samples = num_q_samples))"

        plot();

        target(zs) = exp.(joint_log_density(zs, toy_evidence))
        skillcontour!(target, colour=:red)
        plot_line_equal_skill!(colour=:green)
        variational(zs) = exp.(factorized_gaussian_log_density(params_cur[1],params_cur[2],zs))
        display(skillcontour!(variational, colour=:blue))
    end
    return params_cur
end

```

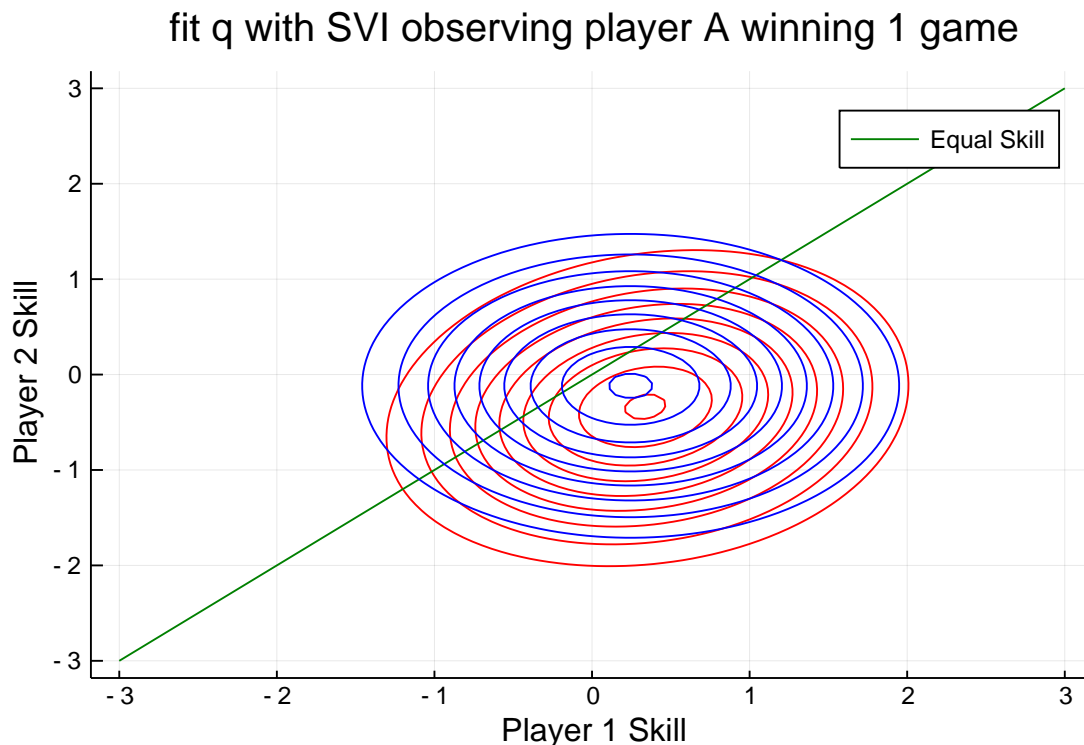
- (d) [2 points] Initialize a variational distribution parameters and optimize them to approximate the joint where we observe player A winning 1 game. Report the final loss. Also plot the optimized variational approximation contours (in blue) and the target distribution (in red) on the same axes.

```

a_win_1 = two_player_toy_games(1, 0)
fit_toy_variational_dist(toy_params_init, a_win_1)
plot!(title="fit q with SVI observing player A winning 1 game"
      xlabel = "Player 1 Skill",
      ylabel = "Player 2 Skill", margin=3mm))

```

Final Loss 0.755657314735221



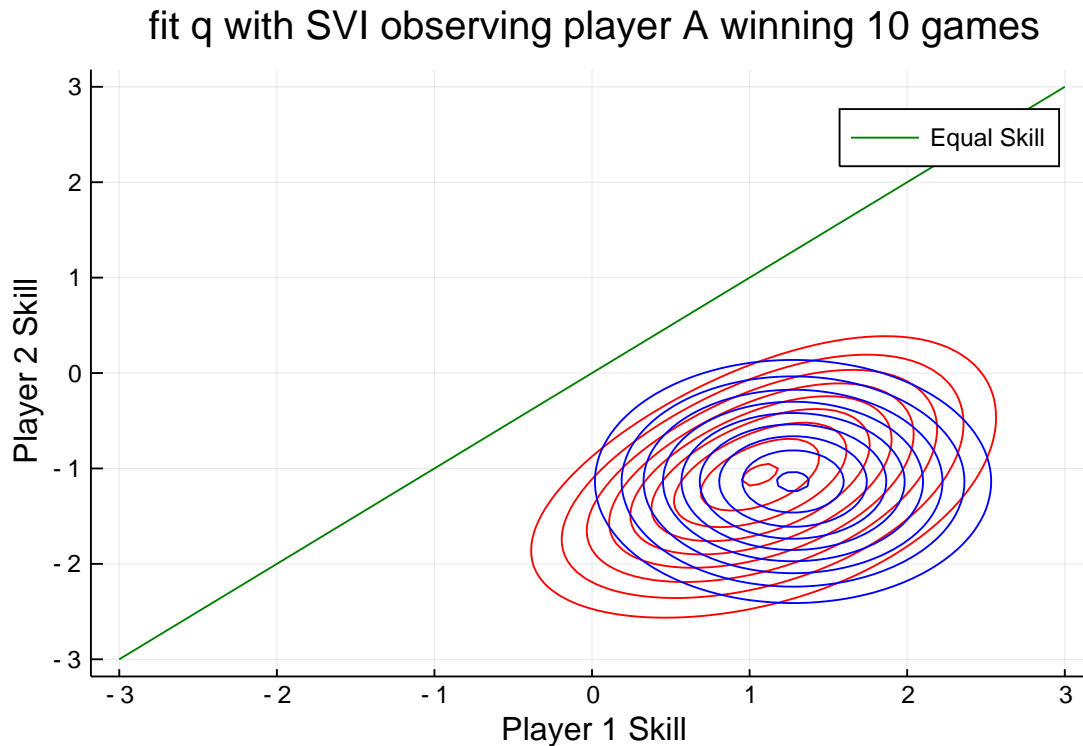
- (e) [2 points] Initialize a variational distribution parameters and optimize them to approximate the joint where we observe player A winning 10 games. Report the final loss. Also plot the optimized variational approximation contours (in blue) and the target distribution (in red) on the same axes.

```

a_win_10 = two_player_toy_games(10, 0)
fit_toy_variational_dist(toy_params_init, a_win_10)
plot!(title="fit q with SVI observing player A winning 10 games",
      xlabel = "Player 1 Skill",
      ylabel = "Player 2 Skill", margin=3mm)

```

Final Loss: 2.766335485285162



- (f) [2 points] Initialize a variational distribution parameters and optimize them to approximate the joint where we observe player A winning 10 games and player B winning 10 games. Report the final loss. Also plot the optimized variational approximation contours (in blue) and the target distribution (in red) on the same axes. For all plots, label both axes.

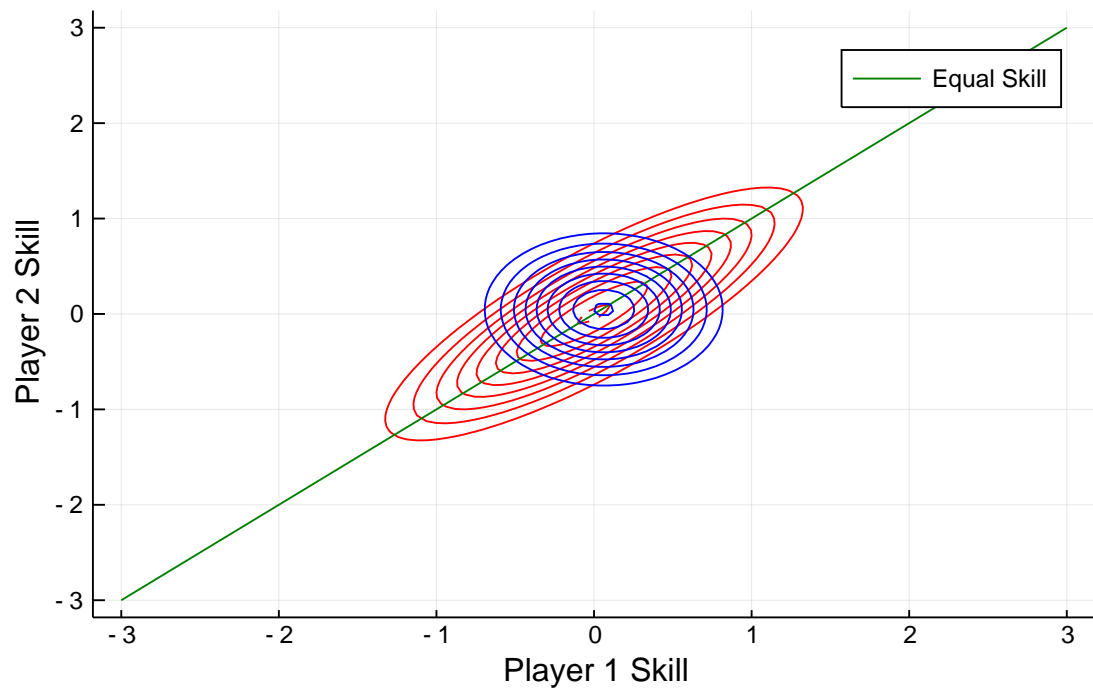
```

split_match_10 = two_player_toy_games(10, 10)
fit_toy_variational_dist(toy_params_init, split_match_10 )
plot!(title="fit q with SVI both winning 10 games",
      xlabel = "Player 1 Skill",
      ylabel = "Player 2 Skill", margin=3mm)

```

Final Loss 15.740345273123456

fit q with SVI both winning 10 games



4 Approximate inference conditioned on real data [24 points]

Load the dataset from `tennis_data.mat` containing two matrices:

- W is a 107 by 1 matrix, whose i 'th entry is the name of player i .
- G is a 1801 by 2 matrix of game outcomes (actually tennis matches), one row per game. The first column contains the indices of the players who won. The second column contains the indices of the player who lost.

Compute the following using your code from the earlier questions in the assignment, but conditioning on the tennis match outcomes:

- (a) [1 point] For any two players i and j , $p(z_i, z_j | \text{all games})$ is always proportional to $p(z_i, z_j | \text{games between } i \text{ and } j)$? In general, are the isocontours of $p(z_i, z_j | \text{all games})$ the same as those of $p(z_i, z_j | \text{games between } i \text{ and } j)$? That is, do the games between other players besides i and j provide information about the skill of players i and j ? A simple yes or no suffices.

Hint: One way to answer this is to draw the graphical model for three players, i , j , and k , and the results of games between all three pairs, and then examine conditional independencies. If you do this, there's no need to include the graphical models in your assignment.

solution: Yes, knowing about other player's skills helps us evaluate player j and i 's skills.

- (b) [5 points] Write a new optimization function `fit_variational_dist` like the one from the previous question except it does not plot anything. Initialize a variational distribution and fit it to the joint distribution with all the observed tennis games from the dataset. Report the final negative ELBO estimate after optimization.

```
using MAT
vars = matread("tennis_data.mat")
player_names = vars["W"]
tennis_games = Int.(vars["G"])
num_players = length(player_names)
print("Loaded data for $num_players players")

function fit_variational_dist(init_params, tennis_games; num_iters=200, lr= 1e-2, num_q_samples = 10)
    params_cur = init_params
    for i in 1:num_iters
        grad_params = gradient((params) -> neg_toy_elbo(params; games = tennis_games, num_samples = num_q_samples), params_cur)[1]
        params_cur = params_cur .- lr .* grad_params
        @info "Loss $(neg_toy_elbo(params_cur; games = tennis_games, num_samples = num_q_samples))"
    end
    return params_cur
end

init_mu = randn(num_players)
init_log_sigma = rand(num_players)
init_params = (init_mu, init_log_sigma)

# Train variational distribution
trained_params = fit_variational_dist(init_params, tennis_games)
means = trained_params[1]
logstd = trained_params[2]

Final Negative ELBO: 1142.7943616545106
```

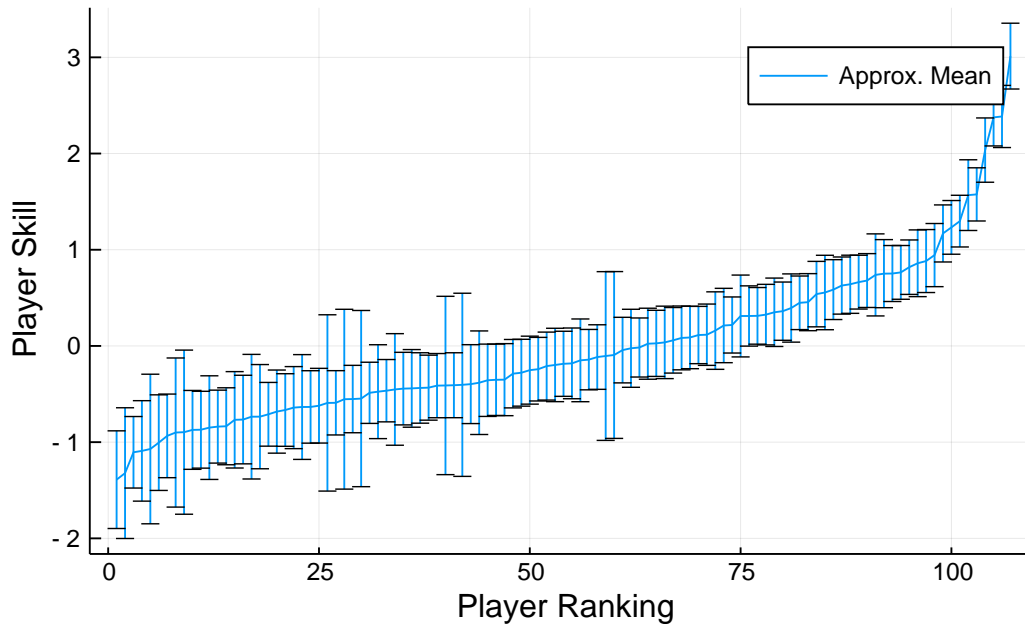
- (c) [2 points] Plot the approximate mean and variance of all players, sorted by skill. For example, in Julia, you can use: `perm = sortperm(means); plot(means[perm], yerror=exp.(logstd[perm]))`. There's no need to include the names of the players.

```
# return players index ordered worst to best
perm = sortperm(means)

# plots player ranking vs player 2 skill, 107th player is top player with highest mean
plot()
plot(means[perm], yerror = exp.(logstd[perm]), title="Approx. Mean & Var of all players",
```

```
xlabel = "Player Ranking",
ylabel = "Player 2 Skill", label = "approx. mean", margin=5mm)
```

Approx. Mean & Var of all players



- (d) [2 points] List the names of the 10 players with the highest mean skill under the variational model.

```
perm = sortperm(means)
ordered_list = sortperm(means, rev=true)
top_ten_players = player_names[ordered_list][1:10]
```

```
"Novak-Djokovic"
"Roger-Federer"
"Rafael-Nadal"
"Andy-Murray"
"David-Ferrer"
"Robin-Soderling"
"Jo-Wilfried-Tsonga"
"Tomas-Berdych"
"Juan-Martin-Del-Potro"
"Richard-Gasquet"
```

- (e) [3 points] Plot the joint posterior over the skills of Roger Federer and Rafael Nadal.

```
fed_i = findall(x->x=="Roger-Federer", player_names)[1][1]
nad_i = findall(x->x=="Rafael-Nadal", player_names)[1][1]
fed_nad_index = [fed_i,nad_i]
skill_m = trained_params[1][fed_nad_index]

games_w_fed = findall(x->x==fed_i, tennis_games)
games_w_nad = findall(x->x==nad_i, tennis_games)
```

```

nlist = []
games = 0
for i in games_w_fed
    games = i[1]
    append!(nlist,games)
end

fed = tennis_games[nlist,:]
fed_nad_games = findall(x->x==1,fed)

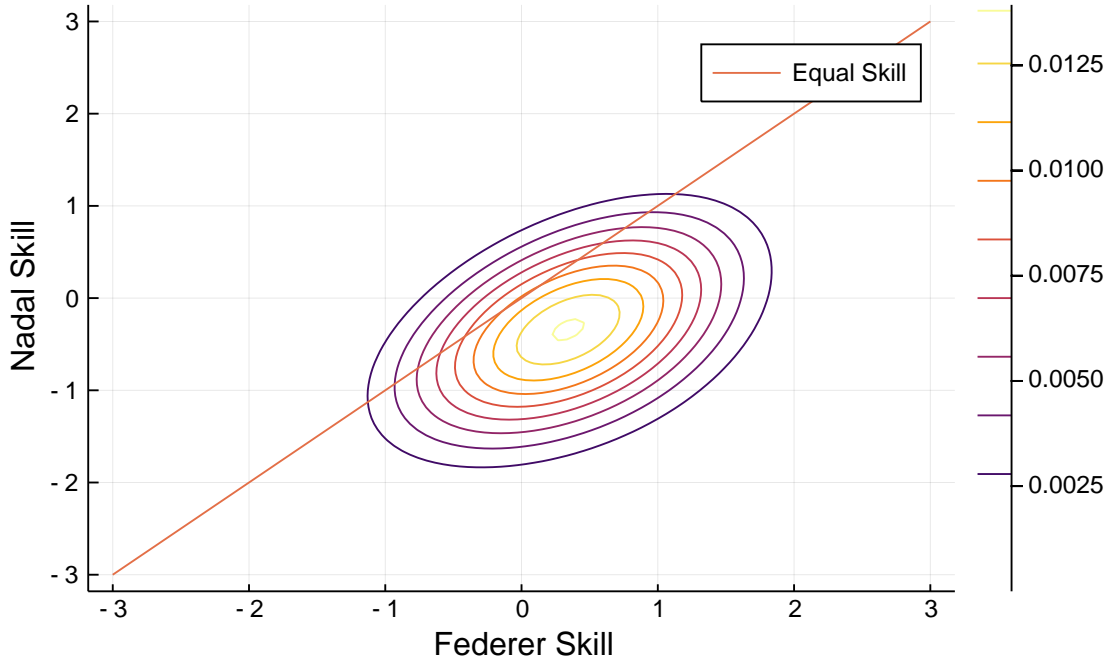
n2list = []
for i in fed_nad_games
    games2 = i[1]
    append!(n2list,games2)
end

FN = fed[n2list,:]
for i in 1:8
    if FN[i] == 5
        FN[i] = 2
    end
end

FN_plot(zs) = exp.(joint_log_density(zs,FN))
plot( title="Federer & Nadal Plot",
      xlabel = "Federer Skill",
      ylabel = "Nadal Skill", margin=5mm)
display(skillcontour!(FN_plot))
plot_line_equal_skill!()

```

Federer & Nadal Plot



(f) [5 points] Derive the exact probability under a factorized Gaussian over two players' skills that one has higher skill than the other, as a function of the two means and variances over their skills.

- Hint 1: Use a linear change of variables $y_A, y_B = z_A - z_B, z_B$. What does the line of equal skill look like after this transformation?
- Hint 2: If $X \sim \mathcal{N}(\mu, \Sigma)$, then $AX \sim \mathcal{N}(A\mu, A^T \Sigma A)$ where A is a linear transformation.
- Hint 3: Marginalization in Gaussians is easy: if $X \sim \mathcal{N}(\mu, \Sigma)$, then the i th element of X has a marginal distribution $X_i \sim \mathcal{N}(\mu_i, \Sigma_{ii})$

Before the transformation, $z_b = z_a$ where z_b is the y axis and z_a is the x axis, we can let z_b be our $f(x)$ and z_a be our x , then the line of best fit from before was $y = x$. After the transformation we have that:

$$z_b = z_a - z_b \Rightarrow 2z_b = z_a \Rightarrow z_b = \frac{z_a}{2} \Rightarrow y = \frac{x}{2} \quad (1)$$

Hence the line of best fit will be flatter.

We want to find A such that:

$$\begin{bmatrix} y_a \\ y_b \end{bmatrix} = A \begin{bmatrix} z_a \\ z_b \end{bmatrix} \quad (2)$$

where $y_a = z_a - z_b$ and $y_b = z_b$.

We find that:

$$A = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} \Rightarrow A \begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix} \Rightarrow \begin{bmatrix} \mu_A - \mu_B \\ \mu_B \end{bmatrix} = \begin{bmatrix} \mu_{y_a} \\ \mu_{y_b} \end{bmatrix} \quad (3)$$

$$\text{cov}(z_a, z_b) = \begin{bmatrix} \sigma_a^2 & 0 \\ 0 & \sigma_b^2 \end{bmatrix} \quad (4)$$

$$(5) \quad A \sum A^T = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \sigma_a^2 & 0 \\ 0 & \sigma_b^2 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix}^T$$

$$= \begin{bmatrix} \sigma_a^2 + \sigma_b^2 & -\sigma_b^2 \\ -\sigma_b^2 & \sigma_b^2 \end{bmatrix} \quad (6)$$

$$\Rightarrow \text{Var}(y_a) = \sigma_a^2 + \sigma_b^2, \text{Var}(y_b) = \sigma_b^2, \text{Cov}(y_a, y_b) = -\sigma_b^2 \quad (7)$$

$$A\mu = A \begin{bmatrix} \mu_{y_a} \\ \mu_{y_b} \end{bmatrix} \quad (8)$$

To derive the probability that A beats B, we arrive at the expression:

$$P(z_a > z_b) = P(z_a - z_b > 0) = P(y_a > 0) = \int_0^\infty f(y_a) dy_a \quad (9)$$

$$y_a \sim \mathcal{N}(A\mu, A \sum A^T) = \mathcal{N}(\mu_A - \mu_B, \sigma_a^2 + \sigma_b^2) \quad (10)$$

[2 points] Compute the probability under your approximate posterior that Roger Federer has higher skill than Rafael Nadal.

Compute this quantity exactly, and then estimate it using simple Monte Carlo with 10000 examples.

```

raf_mu = means[nad_i]
fed_mu = means[fed_i]

raf_sig = exp(logstd[nad_i])
fed_sig = exp(logstd[fed_i])

ya_mu = fed_mu - raf_mu
ya_sig = sqrt(fed_sig^2 + raf_sig^2)

using Distributions: cdf

prob_Fed = 1 - cdf(Normal(ya_mu, ya_sig), 0) (=0.5116349942266001)

function monte_carlo(mu1, sigma1, mu2, sigma2, iters)
    num = 0
    a = 0
    b = 0
    for i in 1:iters
        a = sigma1 .* randn(1) .+ mu1
        b = sigma2 .* randn(1) .+ mu2
        #println(a,b)
        if a > b
            num += 1
        end
    end
    return num/iters
end

monte_carlo(fed_mu, fed_sig, raf_mu, raf_sig, 10000) (=0.5123)

```

[2 points] Compute the probability that Roger Federer is better than the player with the lowest mean skill. Compute this quantity exactly, and then estimate it using simple Monte Carlo with 10000 examples.

```
lowest_mu = means[ordered_list[end]]
lowest_sig = exp(logstd[ordered_list[end]])

ya_mu2 = fed_mu - lowest_mu
ya_sig2 = sqrt(fed_sig^2 + lowest_sig^2)

prob_Fed2 = 1 - cdf(Normal(ya_mu2, ya_sig2), 0) (=0.9999999998254815)

monte_carlo(fed_mu, fed_sig, lowest_mu, lowest_sig, 10000) (=0.9999)
```

[2 points] Imagine that we knew ahead of time that we were examining the skills of top tennis players, and so changed our prior on all players to $\text{Normal}(10, 1)$. Which answers in this section would this change? No need to show your work, just list the letters of the questions whose answers would be different in expectation.

solution: There would be no questions that would be different in expectation.

Collaborated with Rhys Godin